



OpenShift Container Platform 3.11

スケーリングおよびパフォーマンスガイド

OpenShift Container Platform 3.11 スケーリングおよびパフォーマンスガイド

OpenShift Container Platform 3.11 スケーリングおよびパフォーマンスガイド

OpenShift Container Platform 3.11 スケーリングおよびパフォーマンスガイド

法律上の通知

Copyright © 2021 Red Hat, Inc.

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

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

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

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

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

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

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

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

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

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

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

All other trademarks are the property of their respective owners.

概要

実稼働環境でのクラスターのスケールアップおよびパフォーマンスのチューニング

目次

第1章 概要	4
第2章 推奨されるインストール方法	5
2.1. 依存関係の事前インストール	5
2.2. ANSIBLE インストールの最適化	5
2.3. ネットワークの留意事項	6
第3章 ホストの推奨プラクティス	7
3.1. OPENSIFT CONTAINER PLATFORM マスターホストの推奨プラクティス	7
3.2. OPENSIFT CONTAINER PLATFORM ノードホストの推奨プラクティス	8
3.3. OPENSIFT CONTAINER PLATFORM ETCD ホストの推奨プラクティス	9
3.3.1. OpenStack で PCI パススルーを使用した etcd ノードへのストレージ提供	14
3.4. TUNED プロファイルを使用したホストのスケーリング	15
第4章 コンピュートリソースの最適化	17
4.1. オーバーコミット	17
4.2. イメージの留意事項	17
4.2.1. 事前デプロイ済みのイメージを使用した効率の強化	17
4.2.2. イメージの事前プル	18
4.3. RHEL ツールのコンテナイメージを使用したデバッグ	18
4.4. ANSIBLE ベースのヘルスチェックを使用したデバッグ	18
第5章 永続ストレージの最適化	20
5.1. 概要	20
5.2. 一般的なストレージガイドライン	20
5.3. ストレージの推奨事項	21
5.3.1. 特定アプリケーションのストレージの推奨事項	22
5.3.1.1. レジストリー	22
5.3.1.2. スケーリングされたレジストリー	22
5.3.1.3. モニタリング	23
5.3.1.4. ロギング	23
5.3.1.5. アプリケーション	23
5.3.2. 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	24
5.4. グラフドライバーの選択	24
5.4.1. SELinux で OverlayFS または DeviceMapper を使用する利点	28
5.4.2. Overlay と Overlay2 のグラフドライバーの比較	28
第6章 一時ストレージの最適化	30
6.1. 概要	30
6.2. 一般的なストレージガイドライン	30
第7章 ネットワークの最適化	32
7.1. ネットワークパフォーマンスの最適化	32
7.1.1. ネットワークでの MTU の最適化	32
7.2. ネットワークサブネットの設定	33
7.3. IPSEC の最適化	33
第8章 ルーティングの最適化	35
8.1. OPENSIFT CONTAINER PLATFORM HAPROXY ルーターのスケーリング	35
8.1.1. ベースラインのパフォーマンス	35
8.1.2. パフォーマンスの最適化	36
8.1.2.1. 最大接続数の設定	36
8.1.2.2. CPU および割り込みアフィニティー	36

8.1.2.3. スレッド数の増加	37
8.1.2.4. バッファ増加の影響	37
8.1.2.5. HAProxy 再読み込みの最適化	37
第9章 クラスタメトリクスのスケーリング	38
9.1. 概要	38
9.2. OPENSIFT CONTAINER PLATFORM の推奨事項	38
9.3. クラスタメトリクスの容量計画	38
9.4. OPENSIFT CONTAINER PLATFORM メトリクス POD のスケーリング	39
9.4.1. 前提条件	39
9.4.2. Cassandra コンポーネントのスケーリング	39
第10章 CLUSTER MONITORING OPERATOR のスケーリング	41
10.1. 概要	41
10.2. OPENSIFT CONTAINER PLATFORM の推奨事項	41
10.3. CLUSTER MONITORING OPERATOR の容量計画	41
10.3.1. ラボ環境	42
10.3.2. 前提条件	42
第11章 クラスタごとのテスト済み最大数	43
11.1. メジャーリリースについての OPENSIFT CONTAINER PLATFORM のテスト済みクラスタの最大値	43
11.2. OPENSIFT CONTAINER PLATFORM のテスト済みのクラスタの最大値	44
11.2.1. ルートの最大値	45
11.3. OPENSIFT CONTAINER PLATFORM クラスタの最大値をテストする環境および設定	45
11.4. クラスタの最大値に合わせた環境計画	46
11.5. アプリケーション要件に合わせた環境計画	46
第12章 クラスタローダーの使用	48
12.1. クラスタローダーの機能	48
12.2. クラスタローダーのインストール	48
12.3. クラスタローダーの実行	48
12.4. クラスタローダーの設定	48
12.4.1. 設定フィールド	48
12.4.2. クラスタローダー設定ファイルの例	52
12.5. 既知の問題	53
第13章 CPU マネージャーの使用	54
13.1. CPU マネージャーの機能	54
13.2. CPU マネージャーの設定	54
第14章 HUGE PAGE の管理	58
14.1. HUGE PAGE の機能	58
14.2. 前提条件	58
14.3. HUGE PAGE の消費	58
第15章 GLUSTERFS ストレージでの最適化	60
15.1. データベースのコンバインドモードに関するガイド	60
15.2. テスト済みのアプリケーション	60
15.3. サポート表	60
15.4. テスト結果	61

第1章 概要

本ガイドは、OpenShift Container Platform クラスターのパフォーマンスを向上し、OpenShift Container Platform プロダクションスタックの異なるレベルでスケーリングを行う方法についてその手順や例を提供しています。また、本書には、OpenShift Container Platform クラスターのビルド、スケーリング、チューニングの推奨プラクティスが説明されています。

チューニングの留意点は、クラスターの設定により異なり、本書に記載のパフォーマンスに関する推奨事項を実行することで他の部分に影響が及ぶ可能性があるため注意してください。

第2章 推奨されるインストール方法

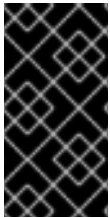
2.1. 依存関係の事前インストール

ノードホストは、ネットワークにアクセスして、**atomic-openshift-***、**iptables**、および **CRI-O** または **Docker** などの RPM 依存関係をインストールします。これらの依存関係を事前にインストールすると、RPM がインストール時にホストごとに何度もアクセスされるのではなく、必要な場合にのみアクセスされるため、より効率的なインストールを実行できます。

また、この方法はセキュリティ上の理由でレジストリーにアクセスできないマシンにも役立ちます。

2.2. ANSIBLE インストールの最適化

OpenShift Container Platform のインストール手法では Ansible を使用します。Ansible は並行して実行する操作に役立ち、迅速かつ効率的なインストールを促進します。ただし、これらの操作はチューニングオプションを追加してさらに強化することができます。利用可能な Ansible 設定オプションの一覧については、「[Ansible の設定](#)」を参照してください。



重要

並行動作は、イメージレジストリーや Red Hat Satellite サーバーなどのコンテンツソースに負荷をかける可能性があります。サーバーのインフラストラクチャー Pod やオペレーティングシステムのパッチを用準備することで、この問題の回避できる可能性があります。

レイテンシーを最小限に抑えたコントロールノード (LAN 速度) からインストーラーを実行します。ワイドエリアネットワーク (WAN) での実行や、ネットワーク接続が途切れる可能性のある環境でのインストールの実行は推奨しません。

Ansible では、RHEL 6.6 以降を使用して OpenSSH のバージョンが **ControlPersist** をサポートすることを確認することや、クラスター内のマシンから **実行せずに**、クラスターと同じ LAN からインストーラーが実行されるようにするなどの、独自の **パフォーマンスやスケーリングに関する指針**が提供されます。

以下は、Ansible で文書化されている推奨事項が組み込まれた、大規模なクラスターのインストールや管理を行うための Ansible の設定例です。

```
# cat /etc/ansible/ansible.cfg
```

出力例

```
# config file for ansible -- http://ansible.com/
# =====
[defaults]
forks = 20 ①
host_key_checking = False
remote_user = root
roles_path = roles/
gathering = smart
fact_caching = jsonfile
fact_caching_connection = $HOME/ansible/facts
fact_caching_timeout = 600
log_path = $HOME/ansible.log
```

```
nocows = 1
callback_whitelist = profile_tasks

[privilege_escalation]
become = False

[ssh_connection]
ssh_args = -o ControlMaster=auto -o ControlPersist=600s -o ServerAliveInterval=60
control_path = %(directory)s/%%h-%%r
pipelining = True ②
timeout = 10
```

- ① フォークは 20 に設定することが理想です (フォークが多くなるとインストールに失敗する可能性があるため)。
- ② パイプラインは、コントロールノードとターゲットノードの間の接続数を減らし、インストーラーのパフォーマンスを向上させます。

2.3. ネットワークの留意事項

ネットワークサブネットの変更はインストール後に実行できますが、これを容易に実行することはできません。サイズを少なく見積もってしまうとクラスターを拡張する際に問題が発生する可能性があるため、ネットワークサブネットのサイズをインストール前に検討することにより作業を大幅に容易にすることができます。

ネットワークサブネットに関する推奨プラクティスは、「[ネットワークの最適化](#)」のトピックを参照してください。

第3章 ホストの推奨プラクティス

3.1. OPENSIFT CONTAINER PLATFORM マスターホストの推奨プラクティス

OpenShift Container Platform インフラストラクチャーで、Pod トラフィックの他に最も使用されるデータパスは OpenShift Container Platform マスターホストと etcd 間のデータパスです。OpenShift Container Platform API サーバー (マスターバイナリーの一部) は、ノードのステータス、ネットワーク設定、シークレットなどについて etcd に確認します。

以下を実行してこのトラフィックパスを最適化します。

- etcd をマスターホストで実行する。デフォルトで、etcd はすべてのマスターホストの静的 Pod で実行されます。
- マスターホスト間でレイテンシーが低く、混雑しない LAN 通信リンクを確保する

OpenShift Container Platform マスターは、CPU 負荷を軽減するためにデシリアライズされたバージョンのリソースを積極的にキャッシュします。ただし、1000 Pod 未満の小規模なクラスターでは、このキャッシュにより、無視できる程度の CPU 負荷を削減するために大量のメモリーが浪費される可能性があります。デフォルトのキャッシュサイズは 50,000 エントリーですが、リソースのサイズによっては 1 から 2 GB メモリーを占有する程度まで拡大する可能性があります。キャッシュのサイズは、`/etc/origin/master/master-config.yaml` で以下の設定を使用して縮小できます。

```
kubernetesMasterConfig:
  apiServerArguments:
    deserialization-cache-size:
      - "1000"
```

API サーバーに送信されるクライアント要求または API 呼び出しの数は、1 秒あたりのクエリー (QPS) 値によって決定され、API サーバーで処理できる同時要求の数は `maxRequestsInFlight` 設定によって決定されます。クライアントが QPS 速度を超える要求を行える数はバースト値によって異なり、この値はバースト性のあるアプリケーションの場合に役立ち、不規則な数の要求の実行を可能にします。要求の応答時間は、大規模または密度の高いクラスターの場合など API サーバーによって処理される同時要求の数が大きい場合にレイテンシーが高くなる可能性があります。Prometheus で `apiserver_request_count` 速度メトリクスを監視し、`maxRequestsInFlight` および **QPS** を適宜調整することが推奨されます。

追加の要求で並列で処理されると API サーバーと etcd IOPS の CPU およびメモリー消費が増大するため、デフォルト値を変更する際にバランスを取ることが必要になります。また、監視 (watch) 以外の要求が過剰にあると 60 秒のタイムアウト後に取り消され、クライアントが再試行する可能性があるため、API サーバーへの負荷が高くなる可能性があることに注意してください。

API サーバーシステムで十分な CPU およびメモリーリソースが利用できることを前提とすると、API サーバー要求のオーバーロードの問題は、上記の要素を考慮し、`*_etc/origin/master/master-config.yaml` で `maxRequestsInFlight`、API qps およびバースト値を引き上げることによって安全に軽減することができます。

```
masterClients:
  openshiftLoopbackClientConnectionOverrides:
    burst: 600
    qps: 300
  servingInfo:
    maxRequestsInFlight: 500
```



注記

上記の `maxRequestsInFlight`、`qps` およびバースト値は OpenShift Container Platform のデフォルトです。`qps` は、要求にかかる時間が1秒未満の場合に `maxRequestsInFlight` 値よりも高くなる場合があります。`maxRequestsInFlight` がゼロに設定されている場合、サーバーが処理できる同時要求の数に制限はありません。

3.2. OPENSIFT CONTAINER PLATFORM ノードホストの推奨プラクティス

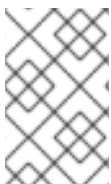
OpenShift Container Platform ノード設定ファイルには、`iptables` 同期期間、SDN ネットワークの Maximum Transmission Unit (MTU)、プロキシモードなどの重要なオプションが含まれます。ノードを設定するには、適切な [ノード設定マップ](#) を変更します。



注記

`node-config.yaml` ファイルを直接編集しないでください。

ノード設定ファイルでは、`kubelet` (node) プロセスに引数を渡すことができます。`kubelet --help` を実行すると、利用可能なオプション一覧を表示できます。



注記

`kubelet` オプションは、OpenShift Container Platform ですべてサポートされておらず、アップストリームの Kubernetes ですべてが使用されている訳ではないので、オプションによってはサポートに制限があります。

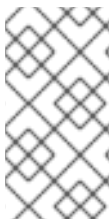


注記

OpenShift Container Platform の各バージョンでサポートされている最大の制限については、「[クラスターの制限](#)」のページを参照してください。

`/etc/origin/node/node-config.yaml` ファイルでは、`pod-s-per-core` および `max-pods` の2つのパラメーターがノードにスケジューリングできる Pod の最大数を制御します。いずれのオプションも使用されている場合、2つの内の小さい方の値でノードの Pod 数が制限されます。これらの値を超えると、以下の状況が発生します。

- OpenShift Container Platform と Docker の両方で CPU 使用率が増加する。
- Pod のスケジューリングの速度が遅くなる。
- メモリー不足のシナリオが生じる可能性がある (ノードのメモリー量によって異なる)。
- 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。



注記

Kubernetes では、単一コンテナを保持する Pod は実際には2つのコンテナを使用します。2つ目のコンテナは実際のコンテナ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のコンテナが実行されていることとなります。

pods-per-core は、ノードのプロセッサ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 プロセッサコアを搭載したノードで **pods-per-core** 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。

```
kubeletArguments:
  pods-per-core:
    - "10"
```



注記

pods-per-core を 0 に設定すると、この制限が無効になります。

max-pods は、ノードのプロパティにかかわらず、ノードが実行できる Pod 数を固定値に設定します。「[クラスタの制限](#)」では **max-pods** のサポートされる最大の値について説明しています。

```
kubeletArguments:
  max-pods:
    - "250"
```

上記の例では、 **pods-per-core** のデフォルト値は **10** であり、 **max-pods** のデフォルト値は **250** です。これは、ノードにあるコア数が 25 以上でない限り、デフォルトでは **pods-per-core** が制限を設定することになります。

OpenShift Container Platform クラスタの推奨制限については、インストールドキュメントの「[サイジングに関する考慮事項](#)」セクションを参照してください。推奨のサイズは、コンテナのステータス更新時の OpenShift Container Platform とコンテナエンジンの連携に対応します。この連携により、大量のログデータの書き込みなどを含め、マスターやコンテナエンジンプロセスへの CPU の負荷がかかります。

kubelet が API サーバーと通信する速度は qps およびバースト値によって異なります。デフォルト値は、各ノードで実行される Pod が制限されている場合には問題なく使用できます。ノードに十分な CPU およびメモリーリソースがある場合、qps およびバースト値は `/etc/origin/node/node-config.yaml` ファイルで微調整できます。

```
kubeletArguments:
  kube-api-qps:
    - "20"
  kube-api-burst:
    - "40"
```



注記

上記の qps およびバースト値は OpenShift Container Platform のデフォルトです。

3.3. OPENSIFT CONTAINER PLATFORM ETCD ホストの推奨プラクティス

etcd は、OpenShift Container Platform が設定に使用するキーと値の分散ストアです。

OpenShift Container Platform のバージョン	etcd のバージョン	ストレージスキーマのバージョン
-------------------------------------	-------------	-----------------

3.3 以前	2.x	v2
3.4 および 3.5	3.x	v2
3.6	3.x	v2 (アップグレード)
3.6	3.x	v3 (新規インストール)
3.7 以降	3.x	v3

etcd 3.x では、クラスターのサイズにかかわらず、CPU、メモリー、ネットワーク、ディスク要件を軽減する、スケーラビリティおよびパフォーマンスの重要な強化機能が導入されました。etcd 3.x は、on-disk etcd データベースの 2 ステップ移行を簡素化する後方互換対応のストレージ API を実装します。移行の目的で、OpenShift Container Platform 3.5 の etcd 3.x で使用されるストレージモードは v2 モードのままとなっています。OpenShift Container Platform 3.6 の時点では、新規インストールで v3 のストレージモードが使用されます。OpenShift Container Platform の以前のバージョンからアップグレードしても、v2 から v3 に自動で **移行されません**。提供されている Playbook を使用して、ドキュメントに記載のプロセスに従い、データを移行する必要があります。

etcd のバージョン 3 では、on-disk etcd データベースの 2 ステップ移行を簡素化する後方互換性対応のストレージ API を実装します。移行の目的で、OpenShift Container Platform 3.5 の etcd 3.x で使用されるストレージモードは v2 モードのままとなっています。OpenShift Container Platform 3.6 の時点で、新規インストールで v3 のストレージモードが使用されます。OpenShift Container Platform 3.7 へのアップグレードプロセスの一環として、必要な場合は etcd ストレージ API を v3 にアップグレードします。バージョン 3.7 以降では、v3 API を使用する必要があります。

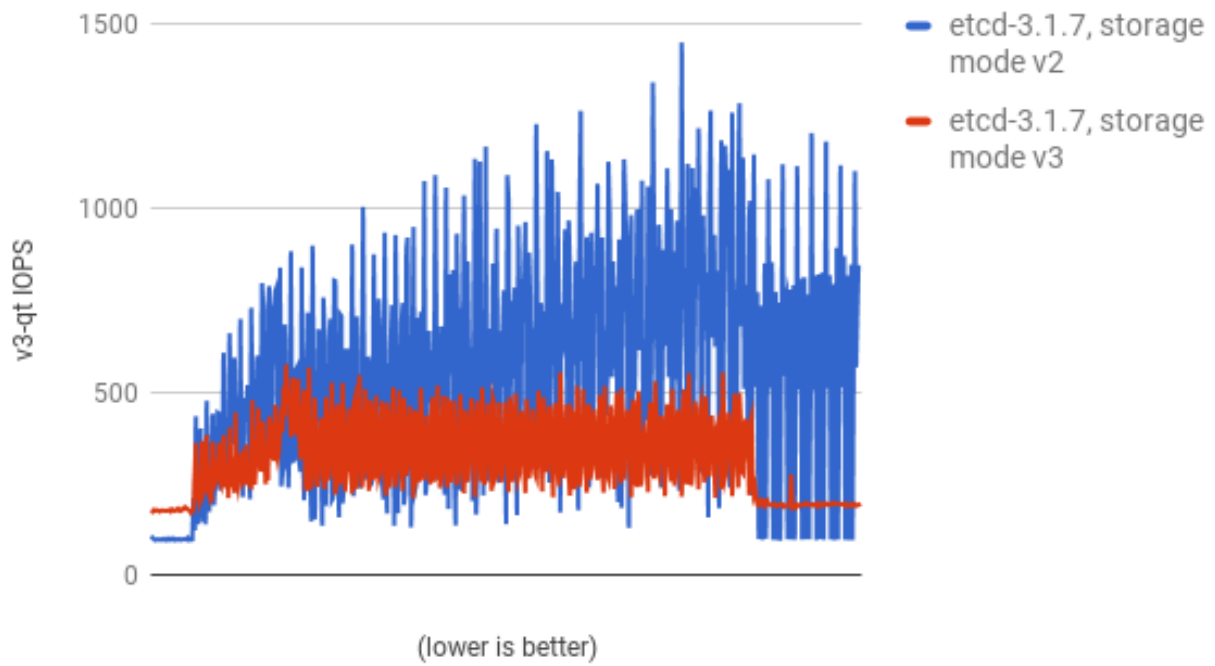
新規インストールでストレージモードを v3 に変更するのに加え、OpenShift Container Platform 3.6 は、全 OpenShift Container Platform タイプに対して強制的に **quorum reads** を実行します。これは、etcd に対するクエリーが古くなったデータを返さないようにするために実行されます。単一ノードの etcd クラスターでは、古くなったデータが入っていても懸念はありませんが、実稼働クラスターで一般的に使用される高可用性の etcd デプロイメントでは、quorum read はクエリーの結果が有効になるようにします。quorum read は、データベース用語の **線形化可能性** (linearizable) と同じです。すべてのクライアントにクラスターが最新の状態に更新されたものが表示され、同じ順番の読み取りおよび書き込みが表示されます。パフォーマンスの向上に関する情報は、etcd 3.1 の [announcement](#) を参照してください。

OpenShift Container Platform は、etcd を使用して Kubernetes 自体が必要な情報以外の追加情報を保存する点を留意することが重要です。たとえば、Kubernetes 以外に OpenShift Container Platform が追加する機能で必要になるので、OpenShift Container Platform は、etcd にイメージ、ビルド、他のコンポーネントの情報を保存します。最終的に、etcd ホストのパフォーマンスやサイジングに関する指針やその他の推奨事項は、Kubernetes とは大幅に異なります。Red Hat は、OpenShift Container Platform のユースケースやパラメーターを考慮に入れて etcd のスケーラビリティやパフォーマンスをテストし、最も正確な推奨事項を提案できるようにしています。

パフォーマンスの向上は、[cluster-loader](#) ユーティリティで 300 ノードの OpenShift Container Platform 3.6 クラスターを使用して、定量化されています。etcd 3.x (ストレージモード v2) と etcd 3.x (ストレージモード v3) を比較すると、以下の図に示されるようにパフォーマンスの向上が明確に確認できます。

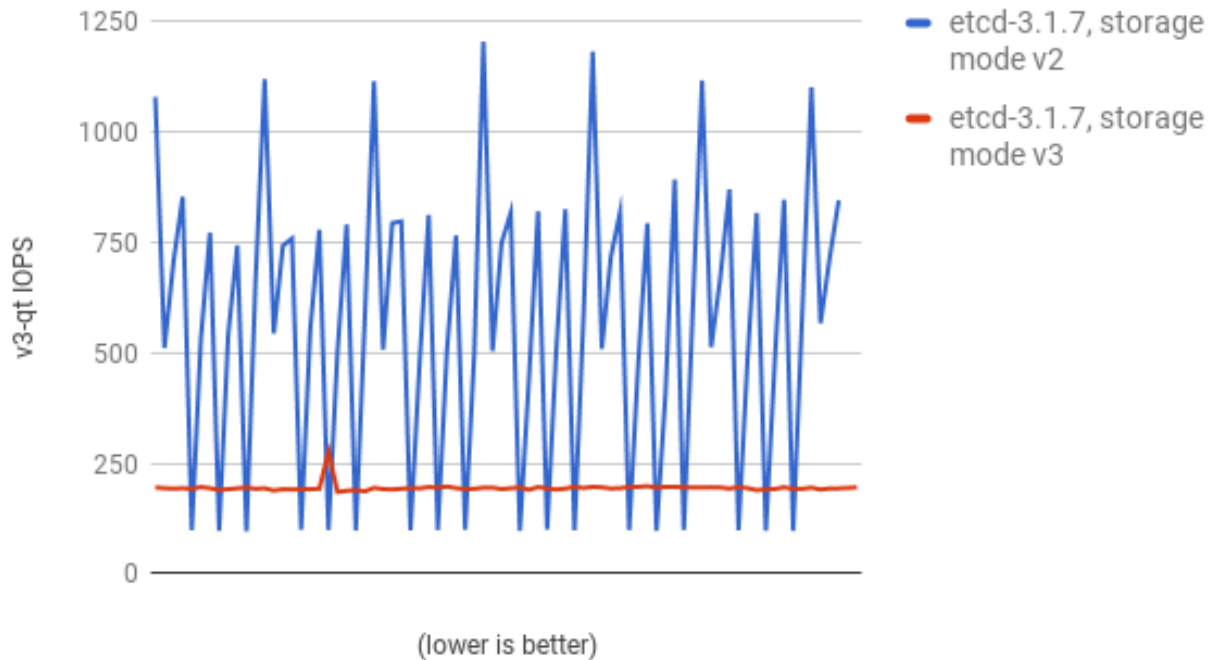
負荷のある状態でのストレージ IOPS が大幅に減少している:

Full run IOPS



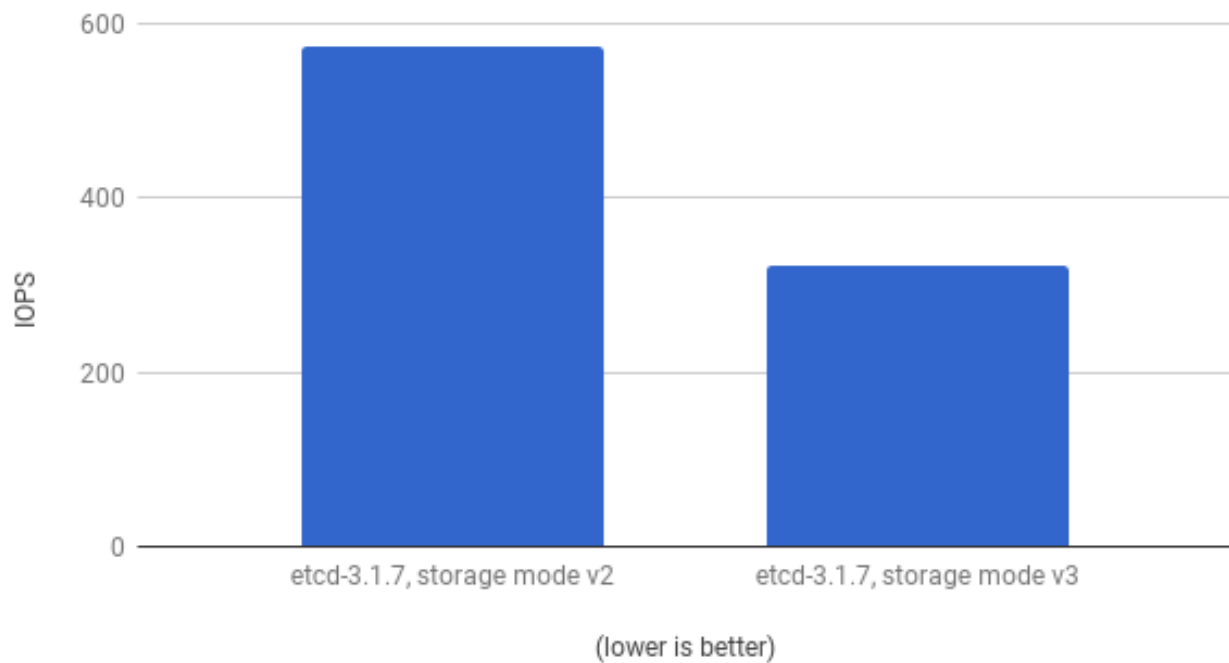
安定した状態でのストレージ IOPS が大幅に減少している:

Steady State IOPS



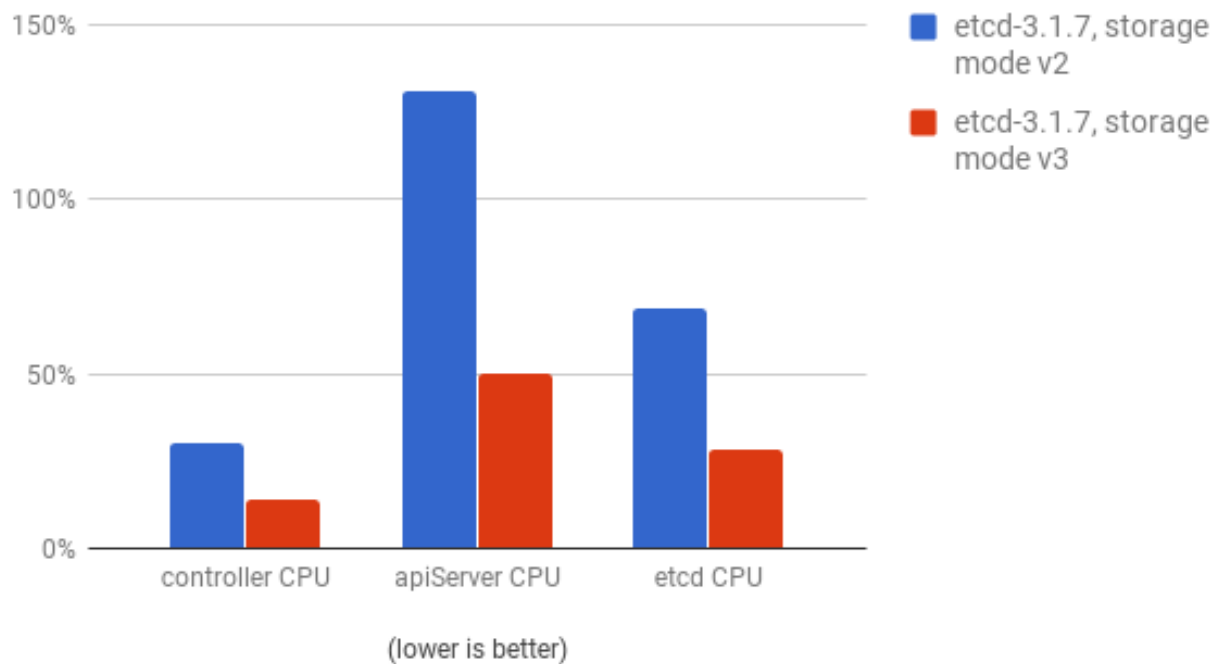
同じ I/O データの表示。両モードでの平均 IOPS

Average Read+Write IOPS



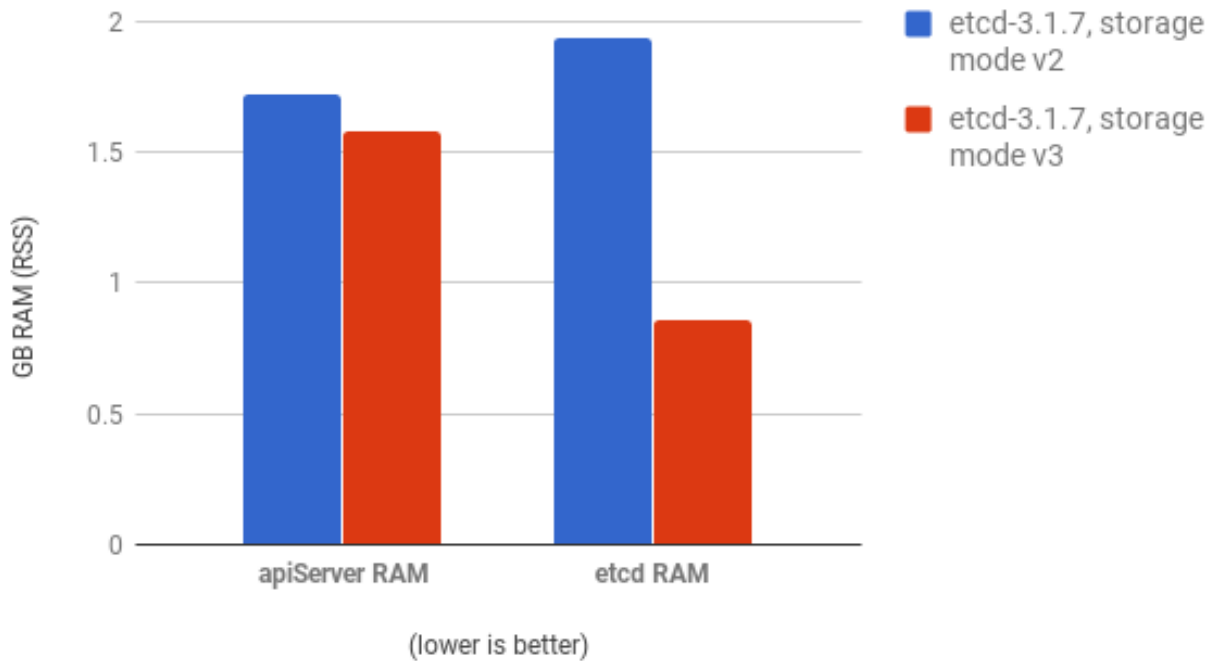
API サーバー (マスター) と etcd プロセスの両方の CPU 使用率が減少している:

CPU Usage



API サーバー (マスター) と etcd プロセスの両方のメモリー使用率も減少している:

Memory Usage (RSS)

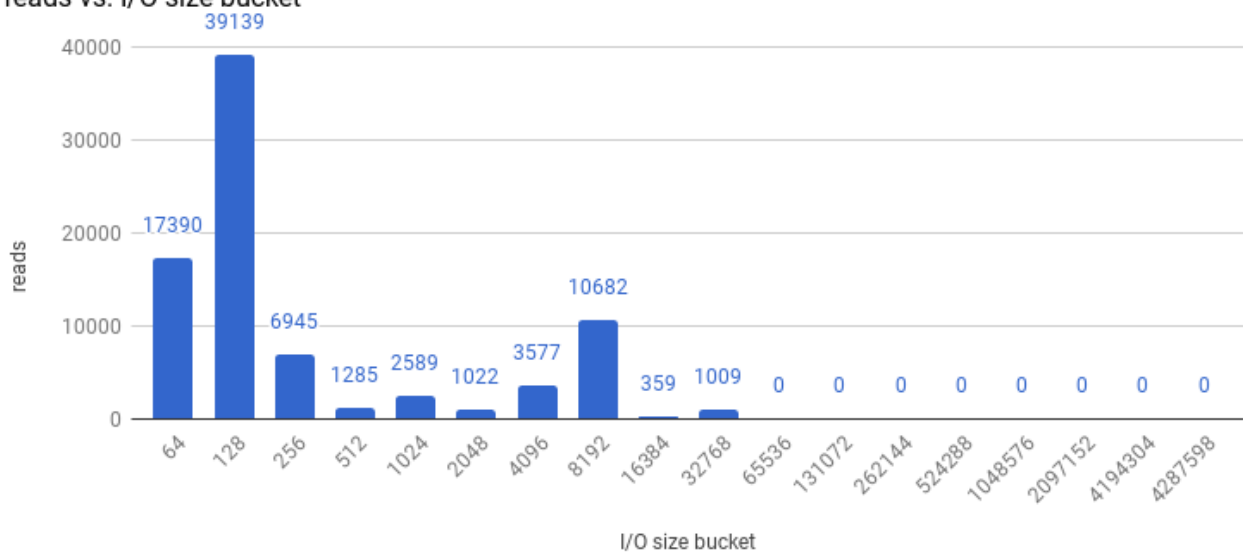


重要

OpenShift Container Platform で etcd をプロファイリングした後に、etcd は少量のストレージインプットおよびアウトプットを頻繁に実行しています。SSD など、少量の読み取り/書き込み操作をすばやく処理するストレージで etcd を使用することを強く推奨します。

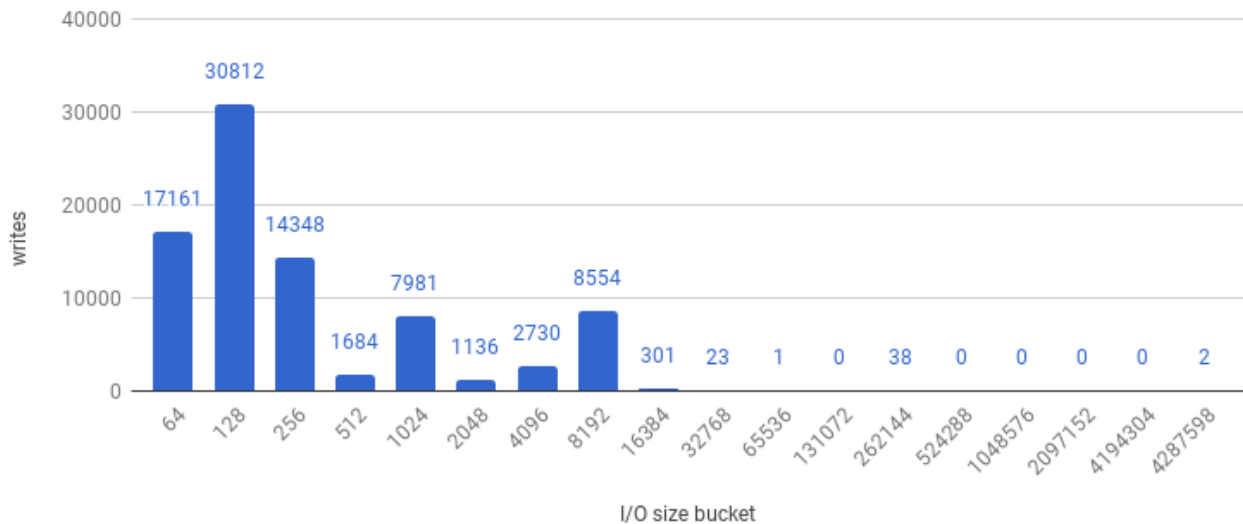
etcd 3.1 の 3 ノードクラスター (quorum reads を有効にしてストレージ v3 モードを使用) で実行したサイズ I/O 操作を確認してみると、読み取りサイズは以下のようになります。

reads vs. I/O size bucket



また、書き込みサイズは以下のようになります。

writes vs. I/O size bucket



注記

etcd プロセスは通常はメモリー集約型であり、マスター/API サーバープロセスは CPU 集約型です。これらは、単一マシンや仮想マシン内で共同設置する上で有効なペアになります。etcd とマスターホスト間の通信を最適化するには、同じホストにマスターと etcd を共同設置するか、専用のネットワークを設定します。

3.3.1. OpenStack で PCI パススルーを使用した etcd ノードへのストレージ提供

大規模な環境で etcd を安定させるために etcd ノードにストレージをすばやく提供するには、NVMe (Non-Volatile Memory express) デバイスを直接 etcd ノードに渡す PCI パススルーを使用します。これを Red Hat OpenStack 11 以降で設定するには、PCI デバイスが存在する OpenStack コンピュートノードで以下を実行してください。

1. Intel Vt-x が BIOS で有効化されているようにします。
2. IOMMU (Input-Output Memory Management Unit) を有効化します。/etc/sysconfig/grub ファイルで、**GRUB_CMDLINX_LINUX** の行末に、引用符で囲って **intel_iommu=on iommu=pt** を追加します。
3. 以下を実行して /etc/grub2.cfg を再生成します。

```
$ grub2-mkconfig -o /etc/grub2.cfg
```

4. システムを再起動します。
5. コントローラーの /etc/nova.conf を以下のように設定します。

```
[filter_scheduler]
```

```
enabled_filters=RetryFilter,AvailabilityZoneFilter,RamFilter,DiskFilter,ComputeFilter,ComputeCapabilitiesFilter,ImagePropertiesFilter,ServerGroupAntiAffinityFilter,ServerGroupAffinityFilter,PciPassthroughFilter
```

```
available_filters=nova.scheduler.filters.all_filters
```

```
[pci]
```

```
alias = { "vendor_id": "144d", "product_id": "a820",
          "device_type": "type-PCI", "name": "nvme" }
```

6. コントローラーで **nova-api** と **nova-scheduler** を再起動します。
7. コンピュートノードの `/etc/nova/nova.conf` で以下のように設定します。

```
[pci]
```

```
passthrough_whitelist = { "address": "0000:06:00.0" }
```

```
alias = { "vendor_id": "144d", "product_id": "a820",
          "device_type": "type-PCI", "name": "nvme" }
```

パススルーする NVMe デバイスの **address**、**vendor_id** および **product_id** の必須値を取得するには、以下を実行します。

```
# lspci -nn | grep devicename
```

8. コンピュートノードで **nova-compute** を再起動します。
9. 実行する OpenStack バージョンで NVMe を使用するように設定し、`etcd` ノードで起動します。

3.4. TUNED プロファイルを使用したホストのスケーリング

Tuned は、Red Hat Enterprise Linux (RHEL) および他の Red Hat 製品で有効な Tuning プロファイルの配信メカニズムです。**Tuned** は、`sysctls`、電源管理、カーネルコマンドラインオプションなどの Linux の設定をカスタマイズして、異なるワークロードのパフォーマンスやスケーラビリティの要件に対応するために、オペレーティングシステムを最適化します。

OpenShift Container Platform は **tuned** デーモンを活用して、**openshift**、**openshift-node** および **openshift-control-plane** と呼ばれる **Tuned** プロファイルを追加します。これらのプロファイルは、カーネルで一般的に発生する垂直スケーリングの上限を安全に増やし、インストール時に自動的にシステムに適用されます。

Tuned プロファイルは、プロファイル間の継承や、プロファイルが仮想環境で使用されるかどうかにより、親プロファイルを選択する親の自動割り当て機能もサポートします。**openshift** プロファイルは、**openshift-node** と **openshift-control-plane** プロファイルの親で、これらの両機能を使用します。OpenShift Container Platform アプリケーションノードとコントロールプレーンノードの両方に関連するチューニングが含まれます。**openshift-node** および **openshift-control-plane** プロファイルは、アプリケーションおよびコントロールプレーンノードにそれぞれ設定されます。

openshift プロファイルがプロファイル階層の親である場合に、OpenShift Container Platform システムに配信されるチューニングは、ベアメタルホスト向けの **throughput-performance** (RHEL のデフォルト) と、RHEL 向けの **virtual-guest** または RHEL Atomic Host ノード向けの **atomic-guest** を組み合わせて作成されます。

お使いのシステムでどの **Tuned** プロファイルが有効になっているかを確認するには以下を実行します。

```
# tuned-adm active
```

出力例

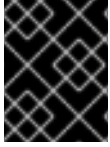
```
Current active profile: openshift-node
```

Tuned に関する詳細情報は、『[Red Hat Enterprise Linux パフォーマンスチューニングガイド](#)』を参照してください。

第4章 コンピュートリソースの最適化

4.1. オーバーコミット

CPU およびメモリーなどのリソースを必要とするクラスターの部分から、このようなリソースにアクセスしやすくなるように、オーバーコミットの手順を使用します。



重要

ハイパーバイザーと Kubernetes 間のスケジュールの衝突による不安定なクラスター動作を防ぐために、ハイパーバイザーレベルではオーバーコミットしないでください。

オーバーコミットすると、別のアプリケーションが必要としているリソースを必要な時にアクセスできなくなってしまうリスクがあり、これにより、結果的にパフォーマンスが低下します。ただし、パフォーマンスが低下する代わりに、密度が高まり、コストが削減されるので、許容できるトレードオフになる場合もあります。たとえば、開発、品質保証 (QA) またはテスト環境でオーバーコミットできても、実稼働環境ではできない場合があります。

OpenShift Container Platform は、コンピュートリソースモデルやクォータシステムでリソース管理を実装します。詳細は、「[OpenShift リソースモデル](#)」を参照してください。

オーバーコミットに関する詳細情報およびストラテジーは、『[クラスター管理ガイド](#)』の「[オーバーコミット](#)」を参照してください。

4.2. イメージの留意事項

4.2.1. 事前デプロイ済みのイメージを使用した効率の強化

効率の強化や全ノードホストでの設定の一貫性の維持、反復タスクの削減を図るため、複数のタスクを組み込んだベースの OpenShift Container Platform イメージを作成できます。これは、事前デプロイ済みのイメージとして知られています。

たとえば、Pod を実行するために、すべてのノードには **ose-pod** イメージが必要なので、各ノードは定期的にコンテナイメージレジストリーに接続して最新のイメージをプルする必要があります。100 のノードが同時にレジストリーに接続し、最新のイメージをプルしようとする問題が発生する場合があります。イメージレジストリーでのリソースの競合や、ネットワーク帯域幅の無駄な使用、および Pod の起動にかかる時間の増加などが生じる可能性があります。

事前にデプロイ済みのイメージをビルドするには、以下を実行します。

- 必要とされるタイプおよびサイズのインスタンスを作成します。
- コンテナ用の永続ボリュームとは別に、専用のストレージデバイスが CRI-O または Docker のローカルイメージやコンテナストレージで利用できるようにします。
- システムを完全に更新すると共に、CRI-O または Docker がインストールされていることを確認します。
- ホストがすべての yum リポジトリーにアクセスできるようにします。
- [シンプロビジョニングされた LVM ストレージ](#)を設定します。

- 一般的に使用するイメージ (rhel7 ベースイメージ) および OpenShift Container Platform インフラストラクチャーコンテナイメージ (**ose-pod**、**ose-deployer** など) を事前にデプロイ済みのイメージに事前に設定します。

OpenStack または AWS で実行できるなど、事前デプロイ済みのイメージが適切なクラスター設定やその他のクラスター設定用に設定されていることを確認します。

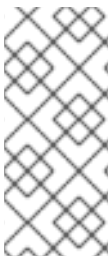
4.2.2. イメージの事前プル

イメージを効率的に生成するには、全ノードホストに、必要とされるコンテナイメージをすべてのノードホストに事前にプルしておきます。これは、最初にイメージをプルする必要がないことを意味するため、サイズが大きくなる可能性のある S2I、メトリクス、ロギングなどのイメージの場合などに、接続速度が遅いことが原因でパフォーマンスが低下することなく、時間を節約できます。

また、この方法はセキュリティ上の理由でレジストリーにアクセスできないマシンにも役立ちます。

または、指定したデフォルトレジストリーではなく、ローカルイメージを使用できます。このためには、以下を実行します。

1. Pod 設定の **imagePullPolicy** パラメーターを **IfNotPresent** または **Never** に設定して、ローカルイメージからプルします。
2. クラスターのすべてのノードで、同じイメージがローカルに保存されていることを確認します。



注記

ノードの設定を制御できる場合は、ローカルレジストリーからプルすることが適切ですが、GCE など、自動的にノードを交換しないクラウドプロバイダーでは確実に機能しない場合があります。Google Container Engine (GKE) で実行している場合は、各ノードに、Google Container Registry の認証情報が設定された **.dockercfg** ファイルが配置されています。

4.3. RHEL ツールのコンテナイメージを使用したデバッグ

Red Hat は **rhel-tools** コンテナイメージを配信します。これは、スケーリングがパフォーマンスの問題のデバッグをサポートするパッケージツールです。このコンテナイメージを使用すると、以下が可能です。

- ベースのディストリビューションからこのサポートコンテナにパッケージを移動して、フットプリントが最小のコンテナホストをデプロイできます。
- 不変のパッケージツリーを含む Red Hat Enterprise Linux 7 Atomic Host 用のデバッグ機能が提供されます。**rhel-tools** には、**tcpdump**、**sosreport**、**git**、**gdb**、**perf** など、多くの一般的なシステム管理ユーティリティーが含まれます。

以下を実行して **rhel-tools** コンテナを使用します。

```
# atomic run rhel7/rhel-tools
```

詳しい情報は、「[RHEL ツールコンテナのドキュメント](#)」を参照してください。

4.4. ANSIBLE ベースのヘルスチェックを使用したデバッグ

追加のヘルスチェックは、OpenShift Container Platform クラスターのインストールおよび管理に使用する [Ansible ベースのツール](#) で利用できます。このヘルスチェックでは、現行の OpenShift Container Platform インストールによくあるデプロイメントの問題を報告します。

これらのチェックは、**ansible-playbook** コマンドの使用 ([クラスターインストール](#) で使用されるのと同じ方式) によるか、または **openshift-ansible** の [コンテナ化されたバージョン](#) として実行できます。**ansible-playbook** 方式については、チェックは **openshift-ansible** RPM パッケージを使って行われます。コンテナ化方式の場合は、**openshift3/ose-ansible** コンテナイメージが [Red Hat Container Registry](#) 経由で配布されます。

利用可能なヘルスチェックや使用例については、『[クラスター管理ガイド](#)』の「[Ansible ベースのヘルスチェック](#)」を参照してください。

第5章 永続ストレージの最適化

5.1. 概要

ストレージを最適化すると、全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。



注記

本ガイドでは、永続ストレージの最適化に重点を置いています。Podの有効期間に使用されるデータ向けのローカルの一時ストレージのオプションは少なくなります。一時ストレージは、一時ストレージのテクノロジープレビューを有効化した場合のみ利用できます。この機能はデフォルトでは無効になっています。詳細情報は、「[一時ストレージの設定](#)」を参照してください。

5.2. 一般的なストレージガイドライン

以下の表では、OpenShift Container Platform で利用可能な永続ストレージ技術を紹介します。

表5.1 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
Block	<ul style="list-style-type: none"> 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バイパス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。 	コンバージドモード/独立モード GlusterFS ^[1] 、iSCSI、ファイバーチャネル、Ceph RBD、OpenStack Cinder、AWS EBS ^[1] 、Dell/EMC Scale.IO、VMware vSphere ポリユーム、GCE 永続ディスク ^[1] 、Azure ディスク
File	<ul style="list-style-type: none"> 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。 	コンバージドモード/独立モード GlusterFS ^[1] 、RHEL NFS、NetApp NFS ^[2] 、Azure File、ベンダー NFS、ベンダー GlusterFS ^[3] 、Azure ファイル、AWS EFS

ストレージタイプ	説明	例
Object	<ul style="list-style-type: none"> 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に組み込む必要があります。 	コンバードモード/インデペンデントモード GlusterFS [1]、Ceph Object Storage (RADOS Gateway)、OpenStack Swift、Aliyun OSS、AWS S3、Google Cloud Storage、Azure Blob Storage、Vendor S3 [3]、Vendor Swift [3]

1. コンバードモード/インデペンデントモード GlusterFS、Ceph RBD、OpenStack Cinder、AWS EBS、Azure Disk、GCE 永続ディスク、および VMware vSphere は、OpenShift Container Platform で永続ボリューム (PV) の動的プロビジョニングをサポートします。
2. 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。
3. Vendor GlusterFS、Vendor S3 および Vendor Swift のサポート機能および設定機能は異なる場合があります。

コンバードモード GlusterFS (ハイパーコンバードまたはクラスターホストのストレージソリューション) またはインデペンデントモード GlusterFS (外部ホストのストレージソリューション) を、OpenShift Container Platform レジストリー、ロギング、モニタリング用のブロック、ファイルおよびオブジェクトストレージに使用できます。

5.3. ストレージの推奨事項

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5.2 設定可能な推奨のストレージ技術

ストレージタイプ	RWO [1]	ROX [2]	RWX [3]	レジストリー	スケーリングされたレジストリー	モニタリング	ロギング	アプリ
Block	Yes	Yes [4]	No	設定可能	設定不可	推奨	推奨	推奨
File	Yes	Yes [4]	Yes	設定可能	設定可能	設定可能 [5]	設定可能 [6]	推奨
Object	Yes	Yes	Yes	推奨	推奨	設定不可	設定不可	設定不可能 [7]

1. ReadWriteOnce

2. ReadOnlyMany
3. ReadWriteMany
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、Azure Disk および Cinder (Block の場合は後者) には該当しません。
5. モニタリングコンポーネントの場合、ReadWriteMany (RWX) アクセスモードのファイルストレージを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は、モニタリングで使用できるように設定されている Persistent Volume Claim (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。logging-es 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、3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。

5.3.1. 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、RHEL NFS サーバーを RHEL でコンテナイメージレジストリー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 サーバー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 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。

5.3.1.1. 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:

- 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。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。NAS ストレージ (オブジェクトストレージインターフェースを使用するのでコンバインドモード/インデペンデントモード GlusterFS 以外) は、実稼働環境のワークロードがある OpenShift Container Platform レジストリークラスターデプロイメントには推奨しません。
- **hostPath** ボリュームは、スケーリングなし/非 HA の OpenShift Container Platform レジストリー用に設定可能ですが、クラスターデプロイメントには推奨しません。

5.3.1.2. スケーリングされたレジストリー

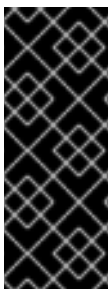
スケーリングされた/高可用性 (HA) の OpenShift Container Platform レジストリーのクラスターデプロイメント:

- 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す。ストレージ技術は、RWX アクセスモードをサポートし、リードアフターライトの一貫性を確保する必要があります。
- 実稼働環境のワークロードを処理するスケーリングされた/HA の OpenShift Container Platform レジストリークラスターのデプロイメントには、ファイルストレージやブロックストレージは推奨しません。
- すべての NAS ストレージ (オブジェクトストレージインターフェースを使用するのでコンバージドモード/インデペンデントモード GlusterFS 以外) は、実稼働環境のワークロードがある OpenShift Container Platform レジストリーのクラスターデプロイメントには推奨しません。

5.3.1.3. モニタリング

OpenShift Container がホストするモニタリ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オブジェクトストレージです。
- ファイルストレージを設定する場合は、POSIX 標準に準拠していることを確認してください。



重要

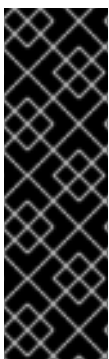
テストの結果、NFS を使用すると修復不能な大規模な破損が発生することが確認されたため、この使用は推奨されません。

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 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。

5.3.1.4. ロギング

OpenShift Container 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オブジェクトストレージです。
- NAS ストレージ (iSCSI からのオブジェクトストレージインターフェースを使用するのでコンバージドモード/インデペンデントモード GlusterFS 以外) は、実稼働環境の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は推奨しません。



重要

テストにより、NFS サーバーを RHEL でコンテナイメージレジストリー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、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 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 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。

5.3.1.5. 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、正常なクラスターをサポートするためにノードに関連付けられていません。
- 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ルアップ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。

5.3.2. 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

- OpenShift Container Platform Internal **etcd**: etcd の信頼性を最も高く保つには、一貫してレイテンシーが最も低くなるストレージ技術が推奨されます。
- 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傾向にあります。

5.4. グラフドライバーの選択

コンテナのランタイムは、イメージとコンテナを DeviceMapper および OverlayFS などのグラフドライバー (プラグ可能なストレージ技術) に保存します。それぞれに長所と短所があります。

サポート内容や使用方法の注意点など、OverlayFS に関する詳しい情報は、『[Red Hat Enterprise Linux \(RHEL\) 7 リリースノート](#)』を参照してください。

表5.3 グラフドライバーの比較

名前	説明	利点	制限
OverlayFS <ul style="list-style-type: none"> • overlay • overlay2 	下層 (親) および上層 (子) のファイルシステムと作業ディレクトリー (子と同じファイルシステム) を組み合わせます。下層のファイルシステムはベースイメージで、新規コンテナを作成すると、差分が含まれる新しい上層ファイルシステムが作成されます。	<ul style="list-style-type: none"> • コンテナの起動および終了時間はデバイスマッパーよりも短くなります。デバイスマッパーと Overlay の起動時間の違いは、1 秒未満です。 • ページキャッシュの共有が可能です。 	POSIX に準拠しません。

名前	説明	利点	制限
デバイスマッパーのシンプロビジョニング	LVM、デバイスマッパー、dm-thinp カーネルモジュールを使用します。ループバックデバイスを削除して、raw パーティション (ファイルシステムなし) に直接移動する点が異なります。	<ul style="list-style-type: none"> ● 中程度の負荷および高密度でパフォーマンス関連のいくつかの利点を測定できます。 ● 容量においてコンテナ別の制限があります (デフォルトは 10GB)。 	<ul style="list-style-type: none"> ● 専用のパーティションが必要で す。 ● Red Hat Enterprise Linux (RHEL) ではデフォルト設定されていません。 ● コンテナおよびイメージはすべて、同じ容量のプールを共有します。プールを破棄または再作成せずにサイズを変更することはできません。
デバイスマッパー loop-lvm	デバイスマッパーのシンプロビジョニングモジュール (dm-thin-pool) を使用して、copy-on-write (CoW) スナップショットを実装します。デバイスマッパーグラフの場所ごとに、2つのブロックデバイス (データとメタデータ用) をベースにシンプールが作成されます。デフォルトでは、これらのブロックデバイスは、自動作成されるスパーファイルのループバックマウントを使用して、自動的に作成されます。	カスタマイズなしですぐに使用できるので、プロトタイプ化や開発の目的で役立ちます。	<ul style="list-style-type: none"> ● Portable Operating System Interface for Unix (POSIX) がすべて機能する訳ではありません (例: O_DIRECT)。最も重要な点として、このモードは実稼働環境のワークロードには対応していません。 ● コンテナおよびイメージはすべて、同じ容量のプールを共有します。プールを破棄または再作成せずにサイズを変更することはできません。

パフォーマンスを強化するために、デバイスマッパーよりも [overlayFS ストレージドライバー](#) を使用することが推奨されます。実稼働環境でデバイスマッパーをすでに使用している場合、コンテナイメージおよびコンテナルートファイルシステムにシンプロビジョニングを使用することを推奨します。そうでない場合、Docker エンジンに overlayfs2 を使用するか、または CRI-O に overlayFS を使用します。

ループデバイスを使用すると、パフォーマンスに影響がある可能性があります。そのまま使用を継続できますが、以下の警告メッセージがログに記録されます。

```
devmapper: Usage of loopback devices is strongly discouraged for production use.
Please use '--storage-opt dm.thinpooldev' or use 'man docker' to refer to
dm.thinpooldev section.
```

ストレージの設定を容易にするには、**docker-storage-setup** ユーティリティーを使用して、設定の詳細の多くを自動化します。

Overlay の場合

1. デバイスドライバーを指定するには、`/etc/sysconfig/docker-storage-setup` ファイルを編集します。

```
STORAGE_DRIVER=overlay2
```



注記

CRI-O を使用している場合は、**STORAGE_DRIVER=overlay** を指定します。

CRI-O の場合、デフォルトの **overlay** ストレージドライバーは **overlay2** の最適化を使用します。

OverlayFS では、**imagefs** を異なる論理ボリュームに配置する必要がある場合、**CONTAINER_ROOT_LV_NAME** および **CONTAINER_ROOT_LV_MOUNT_PATH** を設定する必要があります。**CONTAINER_ROOT_LV_MOUNT_PATH** を設定するには、**CONTAINER_ROOT_LV_NAME** を設定する必要があります。たとえば、**CONTAINER_ROOT_LV_NAME="container-root-lv"** のようになります。詳細は、「[Using the Overlay Graph Driver](#)」を参照してください。

2. Docker ストレージ専用別のディスクドライブがある場合 (例: `/dev/xvdb`) には、以下を `/etc/sysconfig/docker-storage-setup` ファイルに追加します。

```
DEVS=/dev/xvdb
VG=docker_vg
```

3. **docker-storage-setup** サービスを再起動します。

```
# systemctl restart docker-storage-setup
```

4. Docker が **overlay2** を使用していることを確認して、ディスク領域の使用状況をモニタリングするには、**docker info** コマンドを使用します。

```
# docker info | egrep -i 'storage|pool|space|filesystem'
```

出力例

```
Storage Driver: overlay2 1
Backing Filesystem: extfs
```

1 **overlay2** を使用する場合の **docker info** 出力。

Overlay は、Red Hat Enterprise Linux 7.2 の時点で、コンテナランタイムのユースケースについてもサポートされ、起動時間の加速化、ページキャッシュ共有が可能になるため、全体的なメモリー使用率を下げ、高密度化できる可能性があります。

シンプールの場合

1. デバイスドライバーを指定するには、`/etc/sysconfig/docker-storage-setup` ファイルを編集します。

```
STORAGE_DRIVER=devicemapper
```

2. Docker ストレージ専用別のディスクドライブがある場合 (例: `/dev/xvdb`) には、以下を `/etc/sysconfig/docker-storage-setup` ファイルに追加します。

```
DEVS=/dev/xvdb
VG=docker_vg
```

3. `docker-storage-setup` サービスを再起動します。

```
# systemctl restart docker-storage-setup
```

再起動後に、`docker-storage-setup` で、`docker_vg` という名前のボリュームを設定して、シンプールの論理ボリュームを作成します。RHEL でのシンプロビジョニングに関するドキュメントは、『[LVM 管理ガイド](#)』で確認できます。新規作成したボリュームは、`lsblk` コマンドで表示します。

```
# lsblk /dev/xvdb
```

出力例

```
NAME MAJ:MIN RM SIZE RO TYPE MOUNTPOINT
xvdb 202:16 0 20G 0 disk
├─xvdb1 202:17 0 10G 0 part
│   └─docker_vg-docker--pool_tmeta 253:0 0 12M 0 lvm
│       └─docker_vg-docker--pool 253:2 0 6.9G 0 lvm
│           └─docker_vg-docker--pool_tdata 253:1 0 6.9G 0 lvm
│               └─docker_vg-docker--pool 253:2 0 6.9G 0 lvm
```



注記

シンプロビジョニングされたボリュームはマウントされず、ファイルシステムもないので (個別のコンテナには XFS ファイルシステムがない)、`df` の出力には表示されません。

4. Docker が LVM シンプルを使用していることを確認して、ディスク領域の使用状況をモニタリングするには、`docker info` コマンドを使用します。

```
# docker info | egrep -i 'storage|pool|space|filesystem'
```

出力例

```
Storage Driver: devicemapper 1
```

```
Pool Name: docker_vg-docker--pool 2
Pool Blocksize: 524.3 kB
Backing Filesystem: xfs
Data Space Used: 62.39 MB
Data Space Total: 6.434 GB
Data Space Available: 6.372 GB
Metadata Space Used: 40.96 kB
Metadata Space Total: 16.78 MB
Metadata Space Available: 16.74 MB
```

- 1 **devicemapper** を使用する場合は **docker info** 出力。
- 2 `/etc/sysconfig/docker-storage-setup` で指定した **VG** に対応します。

デフォルトでは、シンプールは下層のブロックデバイスの 40% を使用するように設定されています。ストレージを使用していくにつれ、LVM は自動的にプールを最大 100% まで拡張します。**Data Space Total** の値が下層の LVM デバイスのフルサイズに一致しないのは、この理由によります。

開発において、Red Hat ディストリビューションの Docker はループバックマウントが実行されたスパスファイルにデフォルト設定されています。お使いのシステムでループバックモードを使用しているかどうかを確認するには、以下を実行します。

```
# docker info | grep loop0
```

出力例

```
Data file: /dev/loop0
```

5.4.1. SELinux で OverlayFS または DeviceMapper を使用する利点

OverlayFS ファイルシステムの主な利点は、同じノードでイメージを共有するコンテナ間で Linux ページキャッシュが共有される点です。OverlayFS のこの特性により、コンテナの起動時の出入力 (I/O) が減り (数百ミリ秒単位でコンテナの起動時間が短縮)、同様のイメージがノードで実行されている場合にメモリの使用率が減少します。これらはいずれも、(ビルドファームなど) コンテナのチャーンレートを高め、密度を最適化する場合や、イメージの内容に重複が多い環境などの多くの環境で利点があります。

シンプロビジョニングのデバイスがコンテナごとに割り当てられるので、DeviceMapper ではページキャッシュの共有はできません。

5.4.2. Overlay と Overlay2 のグラフィックドライバーの比較

OverlayFS はユニオンファイルシステムの 1 つです。ファイルシステムの上に別のファイルシステムを重ねる (オーバーレイする) ことができます。変更は上層のファイルシステムに記録され、下層のファイルシステムは未変更のままになります。コンテナや DVD-ROM などのファイルシステムイメージを複数のユーザーで共有でき、ベースのイメージは読み取り専用メディアに置かれます。

OverlayFS は、単一の Linux ホストで 2 つのディレクトリーに階層化し、それらを 1 つのディレクトリーとして表示します。これらのディレクトリーは階層と呼ばれ、統合プロセスはユニオンプロセスと呼ばれます。

OverlayFS は、2 つのグラフィックドライバー **overlay** または **overlay2** のいずれかを使用します。Red Hat Enterprise Linux 7.2 の時点では、**overlay** グラフィックドライバーがサポートされるようになりました。Red

Hat Enterprise Linux 7.4 時点で、**overlay2** がサポートされるようになりました。Docker デーモン上の SELinux は、Red Hat Enterprise Linux 7.4 でサポートされるようになりました。サポート内容やご利用のヒントなど、お使いの RHEL バージョンでの OverlayFS の使用に関する情報は、『[Red Hat Enterprise Linux リリースノート](#)』を参照してください。

overlay2 ドライバーは、最大 128 個の下層にある OverlayFS 階層をネイティブでサポートしますが、**overlay** ドライバーは下層の OverlayFS 階層 1 つでしか機能しません。この機能が原因で、**overlay2** ドライバーの方が、**docker build** などの階層関連の Docker コマンドのパフォーマンスが優れており、サポートするファイルシステムで使用する inode が少なくなります。

overlay ドライバーは、下層にある単一の OverlayFS 階層で機能するので、複数の OverlayFS 階層として複数階層のイメージを実装できません。代わりに、各イメージ階層は、`/var/lib/docker/overlay` の配下に独自のディレクトリーとして実装されます。下層にある階層と共有されるデータを参照する場合には、スペース効率が配慮された方法としてハードリンクが使用されます。

Docker は inode の使用において効率が良いので、**overlay** ドライバーではなく、OverlayFS のある **overlay2** ドライバーを使用することが推奨されています。

第6章 一時ストレージの最適化

6.1. 概要



注記

このトピックは、一時ストレージのテクノロジープレビューを有効にした場合にのみ適用されます。この機能は、デフォルトでは無効になっています。この機能を有効にするには、「[configuring for ephemeral storage](#)」を参照してください。



注記

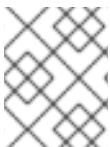
テクノロジープレビューリリースは、Red Hat 製品のサービスレベルアグリーメント (SLA) ではサポートされておらず、機能的に完全でない可能性があり、Red Hat では実稼働環境での使用を推奨しません。テクノロジープレビュー機能は、近々発表予定の製品機能をリリースに先駆けてご提供することにより、お客様に機能性をテストしていただき、開発プロセス中にフィードバックをお寄せいただくことを目的としています。Red Hat テクノロジープレビュー機能のサポート対象範囲に関する詳しい情報は、「[テクノロジープレビュー機能のサポート範囲](#)」を参照してください。

Pod は、一時ファイルの保存などの内部操作に一時ストレージを使用します。この一時ストレージは、個別の Pod の有効期間より長くなることはなく、一時ストレージは Pod 全体で共有できません。

OpenShift Container Platform 3.10 よりも前のバージョンでは、一時ローカルストレージは、コンテナの書き込み可能な階層、logs ディレクトリー、EmptyDir ボリュームを使用して Pod に公開されていました。ローカルストレージのアカウントや分離がないことに関連する問題として、以下の問題が含まれます。

- Pod は利用可能なローカルストレージのサイズを認識しない。
- Pod がローカルストレージを要求しても確実に割り当てられない可能性がある。
- ローカルストレージはベストエフォートのリソースである。
- Pod は他の Pod でローカルストレージが一杯になるとエビクトされる可能性があり、十分なストレージが回収されるまで新しい Pod は許可されない。

一時ストレージは同様に Pod に公開されていますが、Pod の一時ストレージの使用に対する要求や制限を実装する新たな方法が追加されています。



注記

CRI-O をコンテナランタイムとして使用し、ファイルベースロギングをロギングに使用している場合にのみ、コンテナログの管理は該当します。

一時ストレージはシステム内の Pod すべてで共有され、OpenShift Container Platform には管理者およびユーザーが確立した要求や制限を超えるサービスを保証する仕組みはない点を理解することが重要です。たとえば、一時ストレージでは、スループット、秒ごとの I/O 操作またはストレージパフォーマンスについての保証は一切ありません。

6.2. 一般的なストレージガイドライン

ノードのローカルストレージは、プライマリーパーティションとセカンダリーパーティションに分割できます。一時ローカルストレージには、プライマリーパーティションのみが使用できます。プライマリーパーティションでは、`root` とランタイムの2つがサポートされています。

- Root

Root パーティションは、デフォルトで kubelet の root ディレクトリー `/var/lib/kubelet/` と `/var/log/` ディレクトリーを保持します。このパーティションを Pod、オペレーティングシステム、OpenShift Container Platform システムデーモンの間で共有できます。Pod は、EmptyDir ボリューム、コンテナーログ、イメージ階層、コンテナの書き込み可能階層を使用してこのパーティションにアクセスできます。OpenShift Container Platform は、このパーティションの共有アクセスと分離を管理します。

- ランタイム

ランタイムパーティションは、オーバーレイファイルシステムに使用可能なオプションのパーティションです。OpenShift Container Platform は、このパーティションの分離および共有アクセスを特定して提供します。このパーティションには、コンテナイメージ階層と書き込み可能な階層が含まれます。ランタイムパーティションが存在する場合は、**root** パーティションにはイメージ階層も書き込み可能階層も含まれません。

第7章 ネットワークの最適化

7.1. ネットワークパフォーマンスの最適化

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ースカード (NIC) のオフロード、マルチキュー、ethtool の設定を使用してチューニングが可能です。

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。CPU パフォーマンスによっては、このように処理オーバーヘッドが増えることで、従来のオーバーレイ以外のネットワークと比べると、スループットが少なくなり、レイテンシーが増します。

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。

1 Gbps 以上にするには、以下を実行してください。

- **ネイティブのコンテナルーティングの設定**を使用する。このオプションには、ルーターでのルーティングテーブルの更新など、OpenShift SDN の使用時には存在しない重要な操作上の注意事項があります。
- 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、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、システムの CPU から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 サイクルが Pod やアプリケーションで使用できるように解放され、ユーザーはネットワークインフラストラクチャーの帯域幅全体を活用できるようになります。

VXLAN オフロードではレイテンシーは軽減されませんが、CPU 使用率はレイテンシーテストでも減少します。

7.1.1. 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 (ネットワークインターフェースカード (NIC) MTU と、SDN オーバーレイ MTU です)。

NIC MTU は、お使いのネットワークでサポートされる最大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、最大の値を選択するようにしてください。最も低いレイテンシーにおいて最適化するには、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 Ethernet ネットワークでは、この値を 1450 に設定してください。ジャンボフレームの Ethernet ネットワークの場合は、これを 8950 に設定してください。



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。

MTU を設定するには、適切な [ノード設定マップ](#) を編集して、以下のセクションを変更します。

```
networkConfig:
  mtu: 1450 ①
  networkPluginName: "redhat/openshift-ovs-subnet" ②
```

- ① Pod オーバーレイネットワークの Maximum transmission unit (MTU)
- ② `ovs-subnet` プラグインの場合は `redhat/openshift-ovs-subnet` に、`ovs-multitenant` プラグインの場合は `redhat/openshift-ovs-multitenant` に、`ovs-networkpolicy` プラグインの場合は `redhat/openshift-ovs-networkpolicy` に設定します。これは、他の CNI 互換のプラグインにも設定できます。



注記

OpenShift Container Platform SDN を構成するすべてのマスターおよびノードで MTU サイズを変更する必要があります。また、`tun0` インターフェースの MTU サイズはクラスターを構成するすべてのノードで同一である必要があります。

7.2. ネットワークサブネットの設定

OpenShift Container Platform は、Pod とサービスに対して IP アドレス管理を提供します。デフォルト値の許容範囲は以下のとおりです。

- 最大のクラスターサイズは 1024 ノード
- 1024 ノードごとに、/23 を割り当てる (Pod で利用可能な IP は 510 個)
- サービス用の IP アドレスは約 65,536 個

多くの場合、これらのネットワークはデプロイメント後に変更することができません。そのため、拡張についての計画が必要になります。

ネットワークのサイズ変更に関する制限は、「[SDN の設定](#)」に記載されています。

より大規模な環境を計画する場合は、Ansible インベントリーファイルの `[OSE3:vars]` セクションに以下の推奨値を設定するようにしてください。

```
[OSE3:vars]
osm_cluster_network_cidr=10.128.0.0/10
```

これにより、使用可能な IP アドレス 510 個が割り当てられる 8192 ノードが許可されます。

インストールしているソフトウェアバージョンのノード/Pod の上限については、OpenShift Container Platform ドキュメントにあるサポートの制限について参照してください。

7.3. IPSEC の最適化

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。

第8章 ルーティングの最適化

8.1. OPENSIFT CONTAINER PLATFORM HAPROXY ルーターのスケールリング

8.1.1. ベースラインのパフォーマンス

OpenShift Container Platform [ルーター](#) は、宛先が OpenShift Container Platform サービスの外部トラフィックすべてに対する ingress ポイントとなります。

1秒に処理される HTTP 要求について、単一の HAProxy ルーターのパフォーマンスを評価する場合に、パフォーマンスは、以下を含む多くの要因によって変動します。

- 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)

個別の環境でのパフォーマンスは異なりますが、ラボは、サイズが 4 vCPU/16GB RAM のパブリッククラウドインスタンスでテストします。ルート 100 個を処理し、1kB 静的ページに対応するバックエンドで終端される 100 ルートを処理する単一の HAProxy ルーターは、1 秒ごと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	ROUTER_THREADS 設定なし	ROUTER_THREADS=4
なし	23681	24327
edge	14981	22768
passthrough	34358	34331
re-encrypt	13288	24605

HTTP close (keep-alive なし) のシナリオの場合:

暗号化	ROUTER_THREADS 設定なし	ROUTER_THREADS=4
なし	3245	4527

暗号化	ROUTER_THREADS 設定なし	ROUTER_THREADS=4
edge	1910	3043
passthrough	3408	3922
re-encrypt	1333	2239

TLS セッション再開は、暗号化ルートに使用されていました。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、これは多くの場合、プライベートクラウドベースの仮想化にも適用できます。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。ルータ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。

アプリケーションに対してより多くのルートを提供し、ルーティング層の水平スケーリングを図る場合には、[ルーターのシャード](#)を使用する必要があります。

8.1.2. パフォーマンスの最適化

8.1.2.1. 最大接続数の設定

HAProxy のスケーラビリティで最も重要でチューニング可能なパラメーターの 1 つに、**maxconn** パラメーターがあります。このパラメーターは、プロセス別の最大同時接続数を特定の値に設定します。このパラメーターを調整するには、OpenShift Container Platform HAProxy ルーターのデプロイメント設定ファイルにある **ROUTER_MAX_CONNECTIONS** 環境変数を編集してください。



注記

接続にはフロントエンドおよび内部バックエンドが含まれます。これは 2 つの接続としてカウントされます。必ず **ROUTER_MAX_CONNECTIONS** の値を作成しようとしている接続数の 2 倍以上になるように設定してください。

8.1.2.2. CPU および割り込みアフィニティー

OpenShift Container Platform では、HAProxy ルーターは単一のプロセスのとして実行されます。OpenShift Container Platform HAProxy ルーターは通常、周波数が低く、数の多いコアを持つ対称型マルチプロセッシング (SMP) よりも、周波数が高く、数の少ないコアが搭載されたシステムでより優れ

たパフォーマンスを実現します。

HAProxy プロセスを1つの CPU コアに、また別の CPU コアにネットワーク割り込みをピンングすると、ネットワークパフォーマンスが向上する傾向にあります。同じ Non-Uniform Memory Access (NUMA) ノードにプロセスと割り込みを配置すると、共有 L3 キャッシュを確保してメモリアクセスを回避しやすくなります。ただし、このレベルの制御は、パブリッククラウド環境では一般的に不可能です。ベアメタルホストでは、**irqbalance** は、割り込み要求線 (IRQ) があれば、自動的に PCI (peripheral component interconnect) ローカルティイーと NUMA アフィニティーを処理します。クラウド環境では、このレベルの情報は一般的にオペレーティングシステムには提供されません。

CPU ピンングは **taskset** または HAProxy の **cpu-map** パラメーターを使用して実行されます。このディレクティブは、プロセス ID と CPU コア ID の 2 つの引数を取ります。たとえば、HAProxy プロセス 1 を CPU コア 0 にピンングするには、以下の行を HAProxy の設定ファイルの Global セクションに追加します。

```
cpu-map 1 0
```

HAProxy 設定ファイルの変更については、「[カスタマイズされた HAProxy ルーターのデプロイ](#)」を参照してください。

8.1.2.3. スレッド数の増加

HAProxy ルーターは、OpenShift Container Platform のマルチスレッドのサポートと共に提供されます。複数の CPU コアシステムでは、スレッド数が増加すると、ルーター上で SSL を終了する場合などにパフォーマンスが増大する場合があります。

HAProxy ルーターのスレッド数を指定するには、「[Enable HAProxy Threading](#)」および「[Router Environment Variables](#)」を参照してください。

8.1.2.4. バッファ増加の影響

OpenShift Container Platform HAProxy ルーター要求のバッファ設定で、アプリケーションからの受信要求や応答のヘッダーサイズを制限します。HAProxy パラメーター **tune.bufsize** を増やして、より大きいヘッダーを処理し、多くのパブリッククラウドプロバイダーが提供するロードバランサーで許可されるアプリケーションなど、非常に大きい cookie を使用するアプリケーションを機能させることができます。ただし、これにより、多数の接続が開放されている場合など、合計のメモリー使用率に影響があります。非常に多くの接続が開かれている場合には、メモリー使用率は、このチューニング可能なパラメーターの増加とほぼ比例します。

8.1.2.5. HAProxy 再読み込みの最適化

Websocket 接続などの長時間続く接続が、長いクライアント/サーバー HAProxy タイムアウトと短い HAProxy 再読み込み間隔と組み合わせられると、HAProxy プロセスが多数インスタンス化されてしまう可能性があります。これらのプロセスは、HAProxy 設定が再読み込みされる前に開始されていた古い接続を処理する必要があります。これらの数多くのプロセスは、システムに不必要な負荷がかかり、メモリー不足の状態などの問題につながる可能性があるために理想的とは言えません。

この動作に影響を与えるルーターの環境変数は、特に **ROUTER_DEFAULT_TUNNEL_TIMEOUT**、**ROUTER_DEFAULT_CLIENT_TIMEOUT**、**ROUTER_DEFAULT_SERVER_TIMEOUT** および **RELOAD_INTERVAL** などです。

第9章 クラスターメトリクスのスケーリング

9.1. 概要

OpenShift Container Platform は、[Heapster](#) で収集してバックエンドに保存可能なメトリクスを公開します。OpenShift Container Platform の管理者は、1つのユーザーインターフェースでコンテナやコンポーネントメトリクスを表示できます。これらのメトリクスは、スケーリングのタイミングと方法を判断するために、[Horizontal Pod Autoscaler](#) でも使用されます。

以下のトピックでは、メトリクスコンポーネントのスケーリングに関する情報を提供します。



注記

Hawkular および Heapster などのメトリクスコンポーネントの自動スケーリングは OpenShift Container Platform ではサポートされていません。

9.2. OPENSIFT CONTAINER PLATFORM の推奨事項

- 専用の OpenShift Container Platform [インフラストラクチャーノード](#) でメトリクス Pod を実行する
- メトリクスの設定時は永続ストレージを使用する。`USE_PERSISTENT_STORAGE=true` を設定します。
- OpenShift Container Platform メトリクスデプロイメントで `METRICS_RESOLUTION=30` パラメーターを保持する。`METRICS_RESOLUTION` をデフォルト値の `30` よりも小さい値に設定することは推奨していません。Ansible メトリクスのインストール手順を使用する場合は、このパラメーターは `openshift_metrics_resolution` に置き換えてください。
- ホストメトリクス Pod が指定された OpenShift Container Platform ノードを詳しくモニタリングして、ホストシステムの容量不足 (CPU およびメモリー) を早い段階で検出する。このような容量不足により、メトリクス Pod で問題が発生する可能性があります。
- OpenShift Container Platform バージョン 3.7 のテストでは、最大 Pod 数 25,000 のテストケースが OpenShift Container Platform クラスターでモニタリングされました。

9.3. クラスターメトリクスの容量計画

210 および 990 の OpenShift Container Platform ノードで実施したテストでは、10500 および 11000 の Pod がそれぞれモニタリングされ、Cassandra データベースのサイズが、以下の表に記載の速度で増加しました。

表9.1 クラスター内のノード/Podの数に基づく Cassandra データベースのストレージ要件

ノード数	Pod 数	Cassandra ストレージの増加率	1日あたりの Cassandra ストレージの増加量	1週間あたりの Cassandra ストレージの増加量
210	10500	1時間あたり 500 MB	15 GB	75 GB
990	11000	1時間あたり 1GB	30 GB	210 GB

上記の計算では、ストレージ要件が計算値を超過しないようにするための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。

METRICS_DURATION および **METRICS_RESOLUTION** の値がデフォルト (それぞれ 7 日と 15 秒) のままの場合は、上記の値にあるように、安全策として週ごとの Cassandra ストレージのサイズ要件を計画することができます。



警告

OpenShift Container Platform メトリクスはメトリクスデータのデータストアとして Cassandra データベースを使用するので、メトリクス設定のプロセスで **USE_PERSISTENT_STORAGE=true** が設定される場合には、NFS がデフォルトとして設定された状態で PV がネットワークストレージの上部に配置されます。ただし、ネットワークストレージと Cassandra を組み合わせて使用することは推奨していません。

Cassandra データベースをメトリクスデータのデータストアとして使用する場合、推奨事項については、[Cassandra ドキュメント](#)を参照してください。

9.4. OPENSIFT CONTAINER PLATFORM メトリクス POD のスケーリング

メトリクス Pod (Cassandra/Hawkular/Heapster) 1 セットでは、最低 25,000 の Pod をモニタリングできます。

注意

OpenShift Container Platform メトリクス Pod が実行されるノードのシステムの負荷に注意してください。この情報を使用して、OpenShift Container Platform メトリクス Pod の数をスケールアウトし、複数の OpenShift Container Platform ノードに負荷を分散する必要があるかどうかを判断します。OpenShift Container Platform メトリクス heapster Pod のスケーリングは推奨していません。

9.4.1. 前提条件

OpenShift Container Platform メトリクスのデプロイに永続ストレージを使用した場合には、OpenShift Container Platform メトリクスの Cassandra Pod 数をスケールアップする前に、新規 Cassandra Pod が使用されるように、[永続ボリューム \(PV\) を作成する](#)必要があります。ただし、動的にプロビジョニングされる PV を使用して Cassandra がデプロイされた場合には、この手順は必要ありません。

9.4.2. Cassandra コンポーネントのスケーリング

Cassandra ノードは永続ストレージを使用するので、レプリケーションコントローラーでスケールダウンやスケールアップを実行することはできません。

Cassandra クラスターのスケーリングには、**openshift_metrics_cassandra_replicas** 変数を変更して、[デプロイメント](#)を再実行する必要があります。デフォルトでは Cassandra クラスターは単一ノードのクラスターとなっています。

OpenShift Container Platform メトリクスの hawkular Pod を 2 つのレプリカにスケールアップするには、以下を実行します。

```
# oc scale -n openshift-infra --replicas=2 rc hawkular-metrics
```

または、インベントリーファイルを更新して、[デプロイメント](#)を再実行します。



注記

Cassandra クラスターに対して、新規ノードを追加したり、既存のノードを削除した場合は、クラスターに保存したデータの負荷がクラスター全体で再度分散されます。

スケールダウンを実行するには、以下を実行します。

1. コンテナにリモートからアクセスする場合は、削除する Cassandra ノードに対して以下を実行します。

```
$ oc exec -it <hawkular-cassandra-pod> nodetool decommission
```

コンテナにローカルでアクセスする場合には、以下を代わりに実行します。

```
$ oc rsh <hawkular-cassandra-pod> nodetool decommission
```

このコマンドは、データをクラスター全体にコピーするので、実行にしばらく時間がかかります。停止の進捗状況は **nodetool netstats -H** でモニタリングできます。

2. 先のコマンドに成功すると、Cassandra インスタンスの **rc** を **0** にスケールダウンします。

```
# oc scale -n openshift-infra --replicas=0 rc <hawkular-cassandra-rc>
```

これで Cassandra Pod が削除されます。



重要

スケールダウンプロセスが完了し、既存の Cassandra ノードが予想どおりに機能する場合には、この Cassandra インスタンスと対応する Persistent Volume Claim (PVC、永続ボリューム要求) の **rc** も削除できます。PVC を削除すると、この Cassandra インスタンスに関連付けられているデータが完全に削除されるので、スケールダウンが完全かつ正常に完了しなかった場合に、失われたデータを復元することはできません。

第10章 CLUSTER MONITORING OPERATOR のスケーリング

10.1. 概要

OpenShift Container Platform は、[cluster-monitoring-operator](#) によってバックエンドで収集し、保存できるメトリクスを公開します。OpenShift Container Platform 管理者は、Grafana というインターフェースのダッシュボードインターフェースでシステムリソース、コンテナおよびコンポーネントメトリクスを表示できます。

以下のトピックでは、cluster monitoring operator のスケーリングについての情報を提供します。

Prometheus を永続ストレージで使用する場合は、Ansible インベントリーファイルの `openshift_cluster_monitoring_operator_prometheus_storage_enabled` 変数を `true` に設定する必要があります。

10.2. OPENSIFT CONTAINER PLATFORM 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。
- 3つ以上の `openshift-container-storage` ノードを non-volatile memory express (NVMe) ドライブと共に使用します。
- `OpenShift Container Storage (OCS) Block` などの永続的なブロックストレージを使用します。

10.3. CLUSTER MONITORING OPERATOR の容量計画

異なるスケールサイズについて各種のテストが実行されました。以下の表が示すように、Prometheus データベースが拡張されました。



注記

以下の Prometheus ストレージ要件は必ずしも要件を規定するものではありません。ワークロードのアクティビティとリソースの使用状況によっては、クラスターでより高いリソース消費が観察される場合があります。

表10.1 クラスタ内のノード/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

上記の計算では、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` について開発されました。スケールが拡張する場合、Ansible インベントリファイルで、デフォルトが **50Gi** に設定されている `openshift_cluster_monitoring_operator_prometheus_storage_capacity=94Gi` 変数を編集します。



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。

10.3.1. ラボ環境

すべての実験は、OpenStack 環境の OpenShift Container Platform で実行されました。

- インフラストラクチャーノード (VM) - 40 コア、157 GB RAM。
- CNS ノード (VM) - 16 コア、62 GB RAM、NVMe ドライブ。

10.3.2. 前提条件

スケーリングに基づいて、Prometheus データストアの関連する PV サイズを計算し、設定します。デフォルトの Prometheus Pod レプリカ数は 2 であるため、3600 Pod を含む 100 ノードの場合、188 GB が必要になります。

例:

```
195 GB (space per 15 days ) * 2 (pods) = 390 GB free
```

この方程式に基づい

て、`openshift_cluster_monitoring_operator_prometheus_storage_capacity=195Gi` を設定します。

第11章 クラスターごとのテスト済み最大数

OpenShift Container Platform クラスターの計画時に以下のテスト済みのクラスター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、可能な限り規模の大きいクラスターをベースにしています。小規模なクラスターの場合、最大値は比例して低くなります。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。

多くの場合、これらの数値を超えると、全体的にパフォーマンスが低下します。これは、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。

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, Amazon Web Services および Microsoft Azure

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

最大値のタイプ	3.x テスト済みの最大値
ノード数	2,000
Pod 数 ^[1]	150,000
ノードごとの Pod 数	250
コアごとの Pod 数	デフォルト値はありません。
namespaces 数	10,000
ビルド数: パイプラインストラテジー	10,000 (デフォルトの Pod: メモリー 512Mi)
namespace ごとの Pod 数 ^[2]	25,000
サービス数 ^[3]	10,000
namespace ごとのサービス数	5,000
サービスごとのバックエンド数	5,000
namespace ごとのデプロイメント数 ^[2]	2,000

- ここで表示される Pod 数はテスト Pod の数です。Pod の実際の数、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によって異なります。
- 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必要がある、システム内のコントロールループが多数あります。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。この最大値については、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。]

3. 各サービスポートと各サービスのバックエンドに対応するエントリーが含まれます。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。

11.2. OPENSIFT CONTAINER PLATFORM のテスト済みのクラスタの最大値

最大値のタイプ	3.7 テスト済みの最大値	3.9 テスト済みの最大値	3.10 テスト済みの最大値	3.11 テスト済みの最大値
ノード数	2,000	2,000	2,000	2,000
Pod 数 ^[1]	120,000	120,000	150,000	150,000
ノードごとの Pod 数	250	250	250	250
コアごとの Pod 数	デフォルト値は 10 です。	デフォルト値は 10 です。	デフォルト値はありません。	デフォルト値はありません。
namespaces 数	10,000	10,000	10,000	10,000
ビルド数: パイプラインストラテジー	該当なし	10,000 (デフォルトの Pod: メモリー 512Mi)	10,000 (デフォルトの Pod: メモリー 512Mi)	10,000 (デフォルトの Pod: メモリー 512Mi)
namespace ごとの Pod 数 ^[2]	3,000	3,000	3,000	25,000
サービス数 ^[3]	10,000	10,000	10,000	10,000
namespace ごとのサービス数	該当なし	該当なし	5,000	5,000
サービスごとのバックエンド数	5,000	5,000	5,000	5,000
namespace ごとのデプロイメント数 ^[2]	2,000	2,000	2,000	2,000

1. ここで表示される Pod 数はテスト Pod の数です。Pod の実際数は、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によって異なります。
2. 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必要がある、システム内のコントロールループが多数あります。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。この最大値については、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。]

- 各サービスポートおよび各サービスバックエンドには、iptables のバックエンドに対応するエントリーが含まれます。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。

11.2.1. ルートの最大値

OpenShift Container Platform 3.11.53 では、ルーターのテストは Amazon Web Services (AWS) の 3 ノード環境で実行されました。100 の HTTP ルート、具体的には 100 のバックエンド Nginx Pod があり、**keepalive** は **100** に設定されます。結果は以下のようになります。

- ターゲットルートごとに1接続=1秒あたり 24,327 要求
- ターゲットルートごとに 40 接続=1秒あたり 20,729要求
- ターゲットルートごとに 200 接続=1秒あたり 17,253要求

11.3. OPENSIFT CONTAINER PLATFORM クラスターの最大値をテストする環境および設定

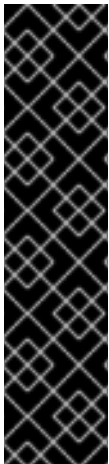
サービスプロバイダーとしてのインフラストラクチャー: OpenStack

ノード	vCPU	RAM(MiB)	ディスクサイズ (GiB)	パススルーディスク	カウント
Master/Etcd ^[1]	16	124672	128	Yes、NVMe	3
Infra ^[2]	40	163584	256	Yes、NVMe	3
クラスター DNS	1	1740	71	No	1
ロードバランサー	4	16128	96	No	1
Container Native Storage ^[3]	16	65280	200	Yes、NVMe	3
Bastion ^[4]	16	65280	200	No	1
ワーカー	2	7936	96	No	2000

- マスター/etcd ノードは、etcd は I/O 集約型でレイテンシーの影響を受けるため、NVMe ディスクでサポートされます。

2. インフラストラクチャーノードはルーター、レジストリー、ロギング、およびモニタリングをホストし、NVMe ディスクでサポートされます。
3. Container Native Storage または Ceph ストレージノードは NVMe ディスクでサポートされません。
4. Bastion ノードは OpenShift Container Platform ネットワークの一部であり、パフォーマンスおよびスケーリングテストのオーケストレーションに使用されます。

11.4. クラスターの最大値に合わせた環境計画



重要

ノードでの物理リソースの過剰なサブスクリプションは、Kubernetes スケジューラーが Pod の配置時に行うリソース保証に影響を与えます。[swap メモリーを無効にする](#)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る可能性があります。

本書で示されている数字は、Red Hat のテスト方法論、セットアップ、設定、およびチューニングに基づくものです。これらの数字は、独自のセットアップや環境によって異なります。

環境の計画時に、ノードに配置できる Pod の数を判断します。

クラスターごとの最大 Pod 数 / ノードごとの予想 Pod 数 = ノード数合計

ノードで適合する Pod 数は、アプリケーション自体により異なります。アプリケーションのメモリー、CPU、ストレージ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 を設定する場合に、ノードごとに最大 250 の Pod があることを前提として、最低でも 9 のノードが必要になります。

$$2200 / 250 = 8.8$$

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。

$$2200 / 20 = 110$$

11.5. アプリケーション要件に合わせた環境計画

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	永続ストレージ
apache	100	500MB	0.5	1GB

Pod タイプ	Pod 数	最大メモリー	CPU コア	永続ストレージ
node.js	200	1GB	1	1GB
postgresql	100	1GB	2	10GB
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。運用上の敏捷性やインスタンスごとのコストなどの要因を考慮する必要があります。

ノードタイプ	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、環境について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。

第12章 クラスターローダーの使用

12.1. クラスターローダーの機能

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。

12.2. クラスターローダーのインストール

クラスターローダーは `atomic-openshift-tests` パッケージに含まれます。これをインストールするには、以下を実行します。

```
$ yum install atomic-openshift-tests
```

インストールが終わると、テスト用の実行ファイル `extended.test` は `/usr/libexec/atomic-openshift/extended.test` に配置されます。

12.3. クラスターローダーの実行

1. **KUBECONFIG** 変数は、管理者 `kubeconfig` の場所に設定します。

```
$ export KUBECONFIG=${KUBECONFIG-$HOME/.kube/config}
```

2. 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。

```
$ cd /usr/libexec/atomic-openshift/
```

```
$ ./extended.test --ginkgo.focus="Load cluster"
```

または `--viper-config` のフラグを追加して、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。

```
$ ./extended.test --ginkgo.focus="Load cluster" --viper-config=config/test 1
```

- 1** この例では、`config/` というサブディレクトリーに `test.yml` という設定ファイルが配置されています。コマンドラインで、ファイルタイプと拡張子はツールが自動的に判断するので、設定ファイルを拡張子なしで実行します。

12.4. クラスターローダーの設定

複数のテンプレートや Pod を含む、namespaces (プロジェクト) を複数作成します。

クラスターローダーの設定ファイルを `config/` サブディレクトリーに配置します。これらの設定ファイルで参照される Pod ファイルとテンプレートファイルは、`content/` サブディレクトリーにあります。

12.4.1. 設定フィールド

表12.1 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 namespaces (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 の配下に、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ています。
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12.2 **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 ConfigMap の名前、値はこの ConfigMap 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。

フィールド	説明
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12.3 pods および templates のフィールド

フィールド	説明
total	このフィールドは使用しません。
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の 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 PodSpec またはテンプレートのいずれかです。
parameters	キーと値のペア。 parameters の配下で、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12.4 tuningset の下のフィールド

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 を識別するサブオブジェクト
templates	テンプレートに適用される tuningsets を識別するサブオブジェクト

表12.5 tuningsets pods または tuningsets templates の下のフィールド

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12.6 tuningsets pods または tuningsets templates, stepping の下のフィールド

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12.7 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。

12.4.2. 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: ./examples/quickstarts/cakephp-mysql.json

    - num: 1
      basename: clusterloader-dancer-mysql
      tuning: default
      ifexists: reuse
      templates:
        - num: 1
          file: ./examples/quickstarts/dancer-mysql.json

    - num: 1
      basename: clusterloader-django-postgresql
      tuning: default
      ifexists: reuse
      templates:
        - num: 1
          file: ./examples/quickstarts/django-postgresql.json

    - num: 1
      basename: clusterloader-nodejs-mongodb
      tuning: default
      ifexists: reuse
      templates:
        - num: 1
          file: ./examples/quickstarts/nodejs-mongodb.json

    - num: 1
      basename: clusterloader-rails-postgresql
      tuning: default
      templates:
        - num: 1
          file: ./examples/quickstarts/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リングしてから、続行します。
- 3 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。
- 4 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。

12.5. 既知の問題

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が **error: unknown parameter name "IDENTIFIER"** エラーで失敗します。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13章 CPU マネージャーの使用

13.1. CPU マネージャーの機能

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考慮されるワークロードに役立ちます。

- できるだけ長い CPU 時間が必要な場合
- プロセッサ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、単一のプロセッサキャッシュを共有することに利点がある場合

13.2. CPU マネージャーの設定

CPU マネージャーを設定するには、以下を実行します。

1. オプションで、ノードにラベルを指定します。

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

2. ターゲットノードで CPU マネージャーのサポートを有効にします。

```
# oc edit configmap <name> -n openshift-node
```

例:

```
# oc edit cm node-config-compute -n openshift-node
```

出力例

```
...
kubeletArguments:
...
feature-gates:
- CPUManager=true
cpu-manager-policy:
- static
cpu-manager-reconcile-period:
- 5s
system-reserved: ①
- cpu=500m
```

```
# systemctl restart atomic-openshift-node
```

- ① **system-reserved** は必須の設定です。この値は、ご使用の環境によって調整が必要になる場合があります。

3. コア1つまたは複数并要求する 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 専用のコアの数になります。

```
# cat cpumanager.yaml
```

出力例

```
apiVersion: v1
kind: Pod
metadata:
  generateName: cpumanager-
spec:
  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"
```

4. Pod を作成します。

```
# oc create -f cpumanager.yaml
```

5. 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-4gdtn
Namespace: test
Node:      perf-node.example.com/172.31.62.105
...
Limits:
  cpu:      1
  memory:  1G
Requests:
  cpu:      1
  memory:   1G
...
QoS Class:   Guaranteed
Node-Selectors: cpumanager=true
               region=primary
```

6. **cgroups** が正しく設定されていることを確認します。一時停止プロセスの PID を取得します。

```
# systemd-cgls -l
```

出力例

```

├─1 /usr/lib/systemd/systemd --system --deserialize 20
├─kubepods.slice
│ └─kubepods-pod0ec1ab8b_e1c4_11e7_bb22_027b30990a24.slice
│   └─docker-
│     └─b24e29bc4021064057f941dc5f3538595c317d294f2c8e448b5e61a29c026d1c.scope
│       └─44216 /pause

```

QoS 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-
pod0ec1ab8b_e1c4_11e7_bb22_027b30990a24.slice/docker-
b24e29bc4021064057f941dc5f3538595c317d294f2c8e448b5e61a29c026d1c.scope
# for i in `ls cpuset.cpus tasks` ; do echo -n "$i "; cat $i ; done

```

出力例

```

cpuset.cpus 2
tasks 44216

```

- 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/44216/status
```

出力例

```
Cpus_allowed_list: 2
```

- 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-burstable.slice/kubepods-burstable-
podbe76ff22_dead_11e7_b99e_027b30990a24.slice/docker-
da621bea7569704fc39f84385a179923309ab9d832f6360cccbff102e73f9557.scope/cpuset.cpus

```

```
0-1,3
```

```
# oc describe node perf-node.example.com
```

出力例

```

...
Capacity:
  cpu: 4
  memory: 16266720Ki
  pods: 40
Allocatable:
  cpu: 3500m
  memory: 16164320Ki
  pods: 40

```

```

---
Namespace           Name           CPU Requests CPU Limits Memory Requests
Memory Limits
-----
test                cpumanager-4gdtm  1 (28%)    1 (28%)    1G (6%)    1G (6%)
test                cpumanager-hczts  1 (28%)    1 (28%)    1G (6%)    1G (6%)
test                cpumanager-r9wrq  1 (28%)    1 (28%)    1G (6%)    1G (6%)
...
Allocated resources:
(Total limits may be over 100 percent, i.e., overcommitted.)
CPU Requests CPU Limits Memory Requests Memory Limits
-----
3 (85%)      3 (85%)    5437500k (32%) 9250M (55%)

```

この仮想マシンには、CPU コアが 4 基あります。**system-reserved** を 500 ミリコアに設定して、**Node Allocatable** の数になるように、ノードの全容量からコアの半分を引きます。

Allocatable CPU が 3500 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 の内 3 つを実行できるという意味になります。1 つのコア全体は、1000 ミリコアに相当します。

4 つ目の Pod をスケジュールしようとする、システムは Pod を受け入れませんが、スケジュールはされません。

```
# oc get pods --all-namespaces |grep test
```

出力例

```

test                cpumanager-4gdtm  1/1    Running    0    8m
test                cpumanager-hczts  1/1    Running    0    8m
test                cpumanager-nb9d5  0/1    Pending    0    8m
test                cpumanager-r9wrq  1/1    Running    0    8m

```

第14章 HUGE PAGE の管理

14.1. HUGE PAGE の機能

メモリーは、ページ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 262,144 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 (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 のサイズが固定されているので、ページサイズを増やすしか、TLB ミスの割合を減らす方法はありません。

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せんが、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割り当て済みの Huge Page を割り当て、消費できます。以下のトピックでは、その方法について説明します。

14.2. 前提条件

1. 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。

14.3. HUGE PAGE の消費

Huge Page は、リソース名の **hugepages-<size>** を使用してコンテナレベルのリソース要件で消費可能です。サイズは、特定のノードでサポートされる最もコンパクトなバイナリー表示 (整数値を使用) に置き換えます。たとえば、ノードが 2048KiB のページサイズをサポートする場合は、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。

```
kind: Pod
metadata:
  generateName: hugepages-volume-
spec:
  containers:
  - securityContext:
    privileged: true
    image: rhel7:latest
    command:
    - sleep
    - inf
    name: example
    volumeMounts:
    - mountPath: /hugepages
```

```

name: hugepage
resources:
  limits:
    hugepages-2Mi: 100Mi ❶
volumes:
- name: hugepage
  emptyDir:
    medium: HugePages

```

- ❶ **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。

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サフィックス [**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。詳しい情報は、「[Huge Page および Transparent Huge Pages](#)」を参照してください。

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。

第15章 GLUSTERFS ストレージでの最適化

15.1. データベースのコンバインドモードに関するガイド

アプリケーションにコンバインドモードを使用する場合には、お使いのワークロードの種類によって、gluster-block と GlusterFS モードを使い分けられるように、このトピックで説明されているガイドとベストプラクティスに従うようにしてください。

15.2. テスト済みのアプリケーション

OpenShift Container Platform 3.10 では、これらの SQL データベースを使用する場合および使用しない場合に関連して広範なテストが行われました。

- PostgreSQL SQL v9.6
- MongoDB noSQL v3.2

これらのデータベースのストレージは、コンバインドモードのストレージクラスターから取得しています。

Postgresql SQL ベンチマークについては、[pgbench](#) がデータベースのベンチマークに使用されていました。MongoDB noSQL ベンチマークについては、YCSB [Yahoo! Cloud Serving Benchmark](#) がベンチマークに使用され、[workloada](#)、[workloadb](#)、[workloadf](#) がテストされました。

15.3. サポート表

表15.1表タイトル: GlusterFS

データベース	ストレージバックエンド: GlusterFS	off にするパフォーマンス変換	on にするパフォーマンス変換
Postgresql SQL	Yes	<ul style="list-style-type: none"> ● performance.start-prefetch ● performance.read-ahead ● performance.write-behind ● performance.read-dir-ahead ● performance.io-cache ● performance.quick-read ● performance.open-behind 	<ul style="list-style-type: none"> ● performance.strict-o-direct

MongoDB noSQL	Yes	<ul style="list-style-type: none"> ● performance.stat-prefetch ● performance.read-ahead ● performance.write-behind ● performance.read-dir-ahead ● performance.io-cache ● performance.quick-read ● performance.open-behind 	<ul style="list-style-type: none"> ● performance.strict-o-direct
---------------	-----	--	---

表15.2 表タイトル: gluster-block

データベース	ストレージバックエンド: gluster-block
Postgresql	Yes
MongoDB	Yes

上述のように GlusterFS のパフォーマンス変換は、コンバージドモードの最新イメージで提供されるデータベースプロファイルにすでに含まれています。

15.4. テスト結果

Postgresql SQL データベースの場合は、GlusterFS と gluster-block のパフォーマンスはほぼ同じ結果となりました。MongoDB noSQL データベースの場合は、gluster-block のパフォーマンスの方が優れていたため、MongoDB noSQL データベースには、gluster-block ベースのストレージを使用してください。