

# AppDirectモード用のVMware ESXiでのDCPMMの設定

## 内容

[概要](#)

[前提条件](#)

[要件](#)

[使用するコンポーネント](#)

[背景説明](#)

[設定](#)

[サービス プロファイルの設定](#)

[ESXiの確認](#)

[仮想マシンNVDIMMの設定](#)

[仮想マシンでの名前空間の構成](#)

[トラブルシューティング](#)

[関連情報](#)

## 概要

このドキュメントでは、ホストマネージドモードでIntel® Optane™ Persistent Memory(PMEM)を使用して、Unified Computing System(UCS)BシリーズサーバでESXiを設定するプロセスについて説明します。

## 前提条件

### 要件

次の項目に関する知識があることが推奨されます。

- UCS B シリーズ
- インテル® Optane™ データ・センター・パーシステント・メモリー・モジュール (DCPMM) の概念
- VMware ESXi および vCenter Server の管理

この設定を行う前に、次の要件が満たされていることを確認します。

- B200/B480 M5仕様ガイドのPMEMガイドラインを参照して[ください](#)。
- CPUが第2世代のIntel® Xeon® スケーラブルプロセッサであることを確認してください。
- PMEM/Dynamic Random Access Memory(DRAM)比は、[KB 67645](#)に基づく要件を満たします。
- ESXiは6.7 U2 + Express Patch 10(ESXi670-201906002)以降です。以前の6.7リリースはサポートされていません。
- UCS Managerとサーバのバージョンは4.0(4)以降です。最新の推奨バージョンについては、

www.software.cisco.com/を[参照してください](#)。

## 使用するコンポーネント

このドキュメントの情報は、次のソフトウェアとハードウェアのバージョンに基づいています。

- UCS B480 M5
- UCS Manager 4.1(2b)

このドキュメントの情報は、特定のラボ環境にあるデバイスに基づいて作成されました。このドキュメントで使用するすべてのデバイスは、初期（デフォルト）設定の状態から起動しています。本稼働中のネットワークでは、各コマンドによって起こる可能性がある影響を十分確認してください。

## 背景説明

App Directモードに設定されたUCSサーバでは、VMware ESXi仮想マシンはOptane DCPMM永続的メモリ不揮発性デュアルインラインメモリモジュール(NVDIMM)にアクセスします。

Intel Optane DCPMMは、IPMCTL管理ユーティリティからUnified Extensible Firmware Interface (UEFI)シェルまたはOSユーティリティを使用して設定できます。このツールは、次のアクションのいくつかを実行するように設計されています。

- モジュールの検出と管理
- モジュールファームウェアのアップデートと設定
- 状態の監視
- 目標、地域、および名前空間のプロビジョニングと設定
- PMEMのデバッグとトラブルシューティング

UCSは、サービスプロファイルに接続された永続的なメモリポリシーを使用して設定できるため、使いやすくなります。

オープンソースの不揮発性デバイス制御(NDCTL)ユーティリティは、LIBNVDIMM Linuxカーネルサブシステムを管理するために使用されます。NDCTLユーティリティを使用すると、システムは設定をプロビジョニングし、OS用の領域および名前空間として実行できます。

ESXiホストに追加された永続メモリは、ホストによって検出され、フォーマットされ、ローカルPMemデータストアとしてマウントされます。ESXiはPMEMを使用するためにVirtual Machine Flying System(VMFS)-Lファイルシステム形式を使用し、ホストごとに1つのローカルPMEMデータストアだけがサポートされます。

PMEMデータストアは、他のデータストアとは異なり、従来のデータストアとしてのタスクをサポートしません。vmxおよびvmware.logファイルを含むVMホームディレクトリをPMEMデータストアに配置することはできません。

PMEMは、次の2種類のモードでVMに提示できます。ダイレクトアクセスモードと仮想ディスクモード。

- ダイレクトアクセスモード  
VMは、NVDIMMの形式でPMEM領域を表示することで、このモードに設定できます。このモ

ードを使用するには、VMオペレーティングシステムがPMem対応である必要があります。NVDIMMはバイトアドレス可能メモリとして動作するため、NVDIMMモジュールに保存されたデータは電源サイクルにわたって保持されます。NVDIMMは、PMEMのフォーマット時にESXiによって作成されたPMemデータストアに自動的に保存されます。

- 仮想ディスクモード

ハードウェアバージョンをサポートするために、VM上に常駐する従来のOSおよびレガシーOSを対象としています。VM OSはPMEM対応である必要はありません。このモードでは、従来のSmall Computer System Interface(SCSI)仮想ディスクを作成し、VM OSで使用できます。

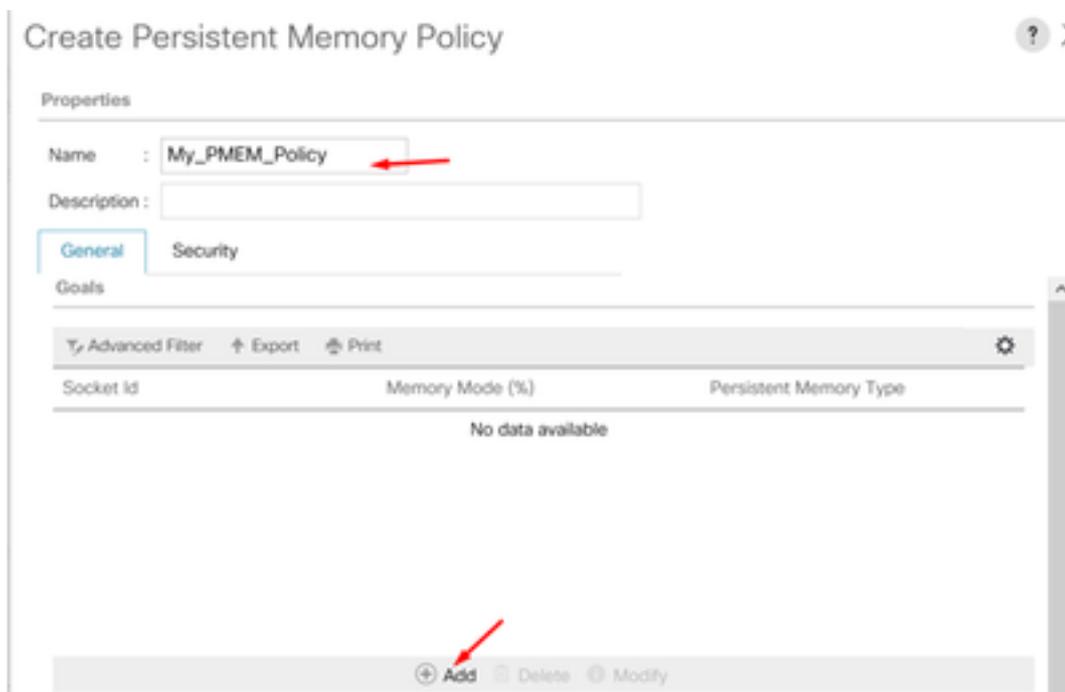
このドキュメントでは、ダイレクトアクセスモードで仮想マシン(VM)を使用するための設定について説明します。

## 設定

この手順では、Intel Optane DCPMMを使用してUCSブレードサーバ上でESXiを設定する方法について説明します。

### サービス プロファイルの設定

1. UCS Manager GUIで、[Servers] > [Persistent Memory Policy]に移動し、図に示すように[Add]をクリックします。



2. 目標を作成し、図に示すようにMemory Modeが0%であることを確認します。

## Create Goal



### Properties

Socket ID :  All Sockets

Memory Mode (%) :

Persistent Memory Type :  App Direct  App Direct Non Interleaved

OK

Cancel

3. PMEMポリシーを目的のサービスプロファイルに追加します。

[Service Profile] > [Policies] > [Persistent Memory Policy]に移動し、作成したポリシーを適用します。

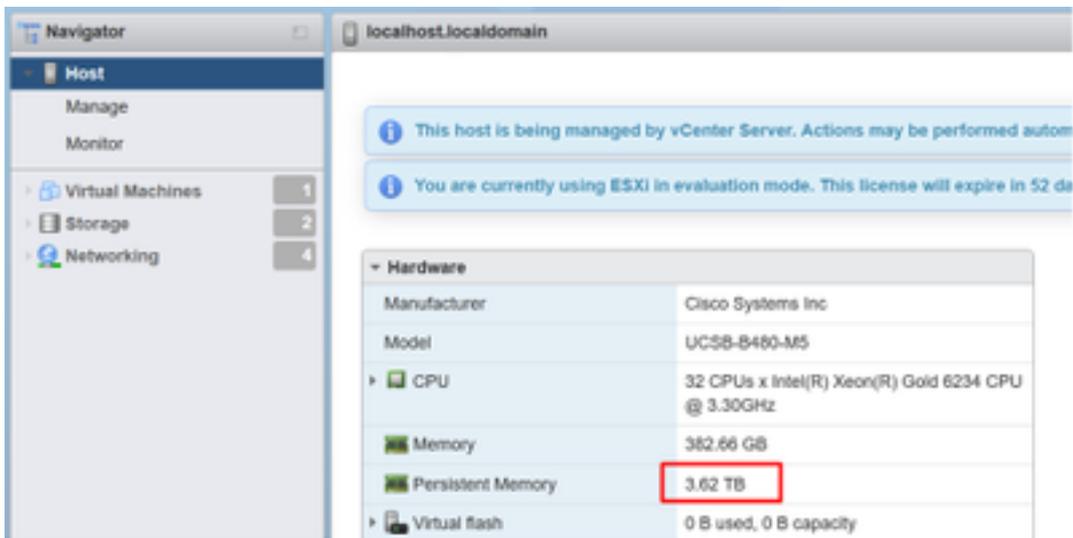
4. 地域の健全性を確認します。

選択した[サーバ] > [インベントリ] > [永続メモリ] > [リージョン]に移動します。タイプ AppDirectが表示されます。このメソッドは、CPUソケットごとに1つの領域を作成します。

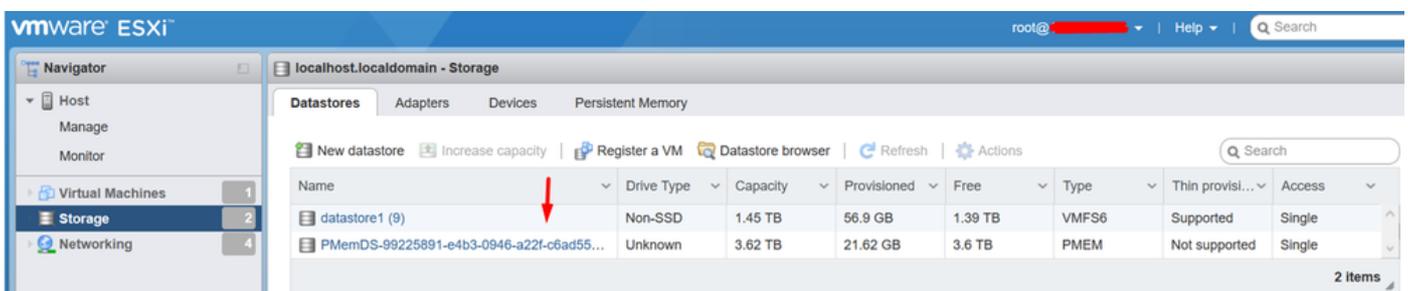
Id	Socket Id	Local DIMM Slot	DIMM Locator Id	Type	Total Capacity (..)	Free Capacity (..)	Health Status
1	Socket 1	Not Applicable	DIMM_A2.DIMM..	AppDirect	928	928	Healthy
2	Socket 2	Not Applicable	DIMM_G2.DIMM..	AppDirect	928	928	Healthy
3	Socket 3	Not Applicable	DIMM_N2.DIMM..	AppDirect	928	928	Healthy
4	Socket 4	Not Applicable	DIMM_U2.DIMM..	AppDirect	928	928	Healthy

## ESXiの確認

1. Webコンソールで、ホストに使用可能な合計PMEMが表示されます。



2. ESXiは、図に示すように、PMEMの総量で構成される特別なデータストアを表示します。



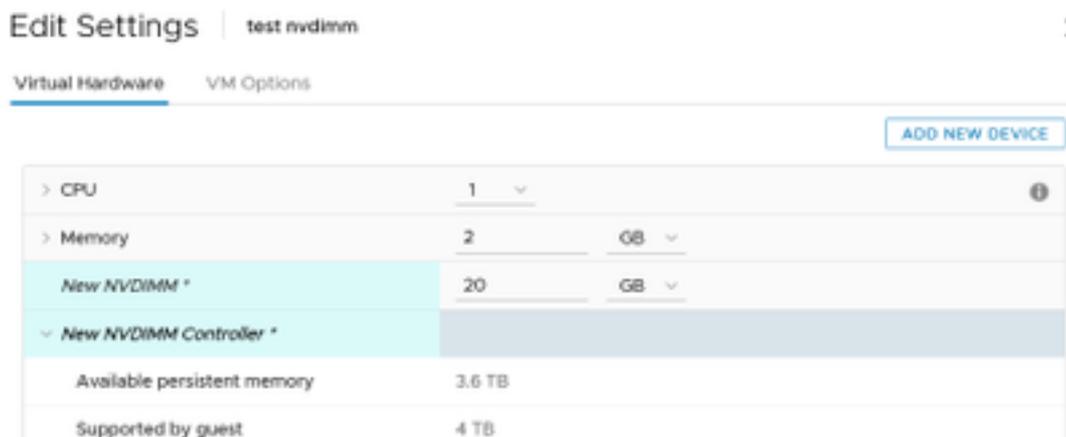
## 仮想マシンNVDIMMの設定

1. ESXiでは、仮想マシンはNVDIMMとしてOptane DCPMM PMEMにアクセスします。NVDIMMを仮想マシンに割り当てるには、vCenterから仮想マシンにアクセスし、[Actions] > [Edit Settings]に移動して、[ADD NEW DEVICE]をクリックして、図に示すように[NVDIMM]を選択します。



注：仮想マシンを作成する場合は、OSの互換性がインテル® Optane™ 永続メモリーをサポートする最小要件バージョンを満たしていることを確認してください。それ以外の場合は、NVDIMMオプションが選択可能な項目に表示されません。

2. 図に示すようにNVDIMMのサイズを設定します。



## 仮想マシンでの名前空間の構成

1. NDCTLユーティリティは、PMEMまたはNVDIMMの管理と設定に使用されます。

この例では、設定にRed Hat 8が使用されています。Microsoftには、永続的なメモリ名前空間の管理用のPowerShellコマンドレットがあります。

Linuxディストリビューションに従って、利用可能なツールを使用してNDCTLユーティリティをダウンロードします

以下に、いくつかの例を示します。

```
# yum install ndctl # zypper install ndctl # apt-get install ndctl
```

2. ESXiがデフォルトで作成したNVDIMM領域と名前空間を確認します。NVDIMMが仮想マシンに割り当てられている場合、スペースが設定と一致していることを確認します。名前空間のモードがrawに設定されていることを確認します。これは、ESXiが名前空間を作成したことを意味します。確認するには、次のコマンドを使用します。

```
# ndctl list -RuN
```

```
admin@localhost:/etc
File Edit View Search Terminal Help
[admin@localhost etc]$ ndctl list -RuN
{
  "regions":[
    {
      "dev":"region0",
      "size":"20.00 GiB (21.47 GB)",
      "available_size":0,
      "max_available_extent":0,
      "type":"pmem",
      "persistence_domain":"unknown",
      "namespaces":[
        {
          "dev":"namespace0.0",
          "mode":"raw",
          "size":"20.00 GiB (21.47 GB)",
          "blockdev":"pmem0"
        }
      ]
    }
  ]
}
```

を選択します。(オプション)名前空間がまだ作成されていない場合は、次のコマンドを

使用してネームスペースを作成できます。

```
# ndctl create-namespace
```

**ndctl create-namespace** コマンドは、デフォルトで **fsdax** モードで新しいネームスペースを作成し、新しい **/dev/pmem([x].[y])** デバイスを作成します。名前空間がすでに作成されている場合は、この手順をスキップできます。

4. PMEM アクセスモードを選択します。設定に使用できるモードは次のとおりです。

- セクタモード :

ストレージを高速ブロックデバイスとして表示します。これは、永続的なメモリを使用できないレガシーアプリケーションに役立ちます。

- Fsdax モード :

永続メモリデバイスが NVDIMM への直接アクセスをサポートできるようにします。ファイルシステムの直接アクセスでは、直接アクセスプログラミングモデルを使用できるようにするために、**fsdax** モードを使用する必要があります。このモードでは、NVDIMM 上にファイルシステムを作成できます。

- Devdax モード :

DAX 文字デバイスを使用して永続メモリへの raw アクセスを提供します。 **devdax** モードを使用するデバイス上でファイルシステムを作成することはできません。

- Raw モード :

このモードには複数の制限があるため、永続メモリの使用は推奨されません。モードを **fsdax** モードに変更するには、コマンドを使用します。

```
ndctl create-namespace -f -e
```

**dev** がすでに作成されている場合は、**dev** 名前空間を使用してモードをフォーマットし、**fsdax** に変更します。

```
admin@localhost:/etc
File Edit View Search Terminal Help
    "size": "20.00 GiB (21.47 GB)",
    "blockdev": "pmem0"
  }
}
}
}
}
[admin@localhost etc]$ ndctl create-namespace -f -e namespace0.0 --mode fsdax
failed to reconfigure namespace: Permission denied
[admin@localhost etc]$ sudo ndctl create-namespace -f -e namespace0.0 --mode fsdax
[sudo] password for admin:
{
  "dev": "namespace0.0",
  "mode": "fsdax",
  "map": "dev",
  "size": "19.69 GiB (21.14 GB)",
  "uuid": "09658ac7-16ea-4c3d-8fbe-e9dae854ddf0",
  "sector_size": 512,
  "blockdev": "pmem0",
  "numa_node": 0
}
[admin@localhost etc]$
```

注：これらのコマンドでは、アカウントにルート権限が必要です。sudoコマンドが必要になる場合があります。

#### 5.ディレクトリとファイルシステムを作成します。

ダイレクトアクセスまたはDAXは、アプリケーションがCPUから（ロードとストアを介して）永続的なメディアに直接アクセスし、従来のI/Oスタックをバイパスできるようにするメカニズムです。DAX対応の永続的メモリファイルシステムには、ext4、XFS、およびWindows NTFSがあります。

作成およびマウントされたXFSファイルシステムの例：

```
sudo mkdir < directory route (e.g./mnt/pmem) > sudo mkfs.xfs < /dev/devicename (e.g. pmem0) >
```

```
admin@localhost:/etc
File Edit View Search Terminal Help
}
[admin@localhost etc]$ mkdir /mnt/pmem
mkdir: cannot create directory '/mnt/pmem': Permission denied
[admin@localhost etc]$ sudo mkdir /mnt/pmem
[admin@localhost etc]$ sudo mkfs.xfs /dev/pmem0
meta-data=/dev/pmem0          isize=512    agcount=4, agsize=1290112 blks
=                               sectsz=4096 attr=2, projid32bit=1
=                               crc=1      finobt=1, sparse=1, rmapbt=0
=                               reflink=1
data      =                   bsize=4096 blocks=5160448, imaxpct=25
=                               sunit=0    swidth=0 blks
naming    =version 2          bsize=4096 ascii-ci=0, ftype=1
log       =internal log     bsize=4096 blocks=2560, version=2
=                               sectsz=4096 sunit=1 blks, lazy-count=1
realtime  =none              extsz=4096 blocks=0, rtextents=0
[admin@localhost etc]$
```

#### 6.ファイルシステムをマウントし、正常に実行されたことを確認します。

```
sudo mount
```

```
admin@localhost:/etc
File Edit View Search Terminal Help
[admin@localhost etc]$ sudo mount /dev/pmem0 /mnt/pmem/
[admin@localhost etc]$ // verify the mount was successful
bash: //: Is a directory
[admin@localhost etc]$ df -h /mnt/pmem/
Filesystem      Size  Used Avail Use% Mounted on
/dev/pmem0      20G  173M   20G   1% /mnt/pmem
[admin@localhost etc]$
```

VMはPMEMを使用する準備ができました。

## トラブルシューティング

エラーが見つかった場合は、`-o dax`マウントオプションを使用してこのDAX対応ファイルシステムをマウントすることをお勧めします。

```
[admin@localhost etc]$ sudo mount -o dax /dev/pmem0 /mnt/pmem/
mount: /mnt/pmem: wrong fs type, bad option, bad superblock on /dev/pmem0, missing codepage or helper program, or other error.
```

整合性を確保するために、ファイルシステムの修復が実行されます。

```
[admin@localhost etc]$ sudo xfs_repair /dev/pmem0
[sudo] password for admin:
Phase 1 - find and verify superblock...
Phase 2 - using internal log
- zero log...
- scan filesystem freespace and inode maps...
- found root inode chunk
Phase 3 - for each AG...
- scan and clear agi unlinked lists...
- process known inodes and perform inode discovery...
- agno = 0
- agno = 1
- agno = 2
- agno = 3
- process newly discovered inodes...
Phase 4 - check for duplicate blocks...
- setting up duplicate extent list...
- check for inodes claiming duplicate blocks...
- agno = 0
- agno = 1
- agno = 2
- agno = 3
Phase 5 - rebuild AG headers and trees...
- reset superblock...
Phase 6 - check inode connectivity...
- resetting contents of realtime bitmap and summary inodes
- traversing filesystem ...
- traversal finished ...
- moving disconnected inodes to lost+found ...
Phase 7 - verify and correct link counts...
done
[admin@localhost etc]$
```

回避策として、`-o dax`オプションを使用せずにマウントをマウントできます。

注：xfsprogsバージョン5.1では、デフォルトでは`reflink`オプションを有効にしてXFSファイルシステムを作成することになります。以前は、デフォルトで無効になっていました。`reflink`と`dax`のオプションは相互に排他的で、マウントが失敗します。

「DAXとreflinkを一緒に使用することはできません！」 dmesgでは、mountコマンドが失敗するとエラーが表示されます。

```
admin@localhost:/etc
File Edit View Search Terminal Help
log      =internal log          bsize=4096   blocks=2560, version=2
         =                    sectsz=4096  sunit=1 blks, lazy-count=1
realtime =none              extsz=4096   blocks=0, rtextents=0
[admin@localhost etc]$ mount -o dax /dev/pmem0 /mnt/pmem
mount: only root can use "--options" option
[admin@localhost etc]$ sudo mount -o dax /dev/pmem0 /mnt/pmem/
mount: /mnt/pmem: wrong fs type, bad option, bad superblock on /dev/pmem0, missing
codepage or helper program, or other error.
[admin@localhost etc]$ dmesg -T | tail
[mar nov 10 00:12:18 2020] VFS: busy inodes on changed media or resized disk sr0
[mar nov 10 00:12:22 2020] ISO 9660 Extensions: Microsoft Joliet Level 3
[mar nov 10 00:12:22 2020] ISO 9660 Extensions: RRIP_1991A
[mar nov 10 01:47:35 2020] pmem0: detected capacity change from 0 to 21137195008
[mar nov 10 01:51:19 2020] XFS (pmem0): DAX enabled. Warning: EXPERIMENTAL, use
at your own risk
[mar nov 10 01:51:19 2020] XFS (pmem0): DAX and reflink cannot be used together!
[mar nov 10 01:53:06 2020] XFS (pmem0): DAX enabled. Warning: EXPERIMENTAL, use
at your own risk
[mar nov 10 01:53:06 2020] XFS (pmem0): DAX and reflink cannot be used together!
[mar nov 10 01:59:29 2020] XFS (pmem0): DAX enabled. Warning: EXPERIMENTAL, use
at your own risk
[mar nov 10 01:59:29 2020] XFS (pmem0): DAX and reflink cannot be used together!
[admin@localhost etc]$
```

回避策として、-o daxオプションを削除します。

```
admin@localhost:/etc
File Edit View Search Terminal Help
[admin@localhost etc]$ sudo mount /dev/pmem0 /mnt/pmem/
[admin@localhost etc]$ // verify the mount was successful
bash: //: Is a directory
[admin@localhost etc]$ df -h /mnt/pmem/
Filesystem      Size  Used Avail Use% Mounted on
/dev/pmem0      20G  173M   20G   1% /mnt/pmem
[admin@localhost etc]$
```

ext4 FSでマウントします。

EXT4ファイルシステムはreflink機能を実装していないがDAXをサポートしているため、代替として使用できます。

```
[admin@localhost etc]$ sudo mkfs.ext4 /dev/pmem0
mke2fs 1.44.3 (10-July-2018)
/dev/pmem0 contains a xfs file system
Proceed anyway? (y,N) y
Creating filesystem with 5160448 4k blocks and 1291808 inodes
Filesystem UUID: 164c6d57-0462-45a0-9b94-703719272816
Superblock backups stored on blocks:
    32768, 98304, 163840, 229376, 294912, 819200, 884736, 1605632, 2654208,
    4096000

Allocating group tables: done
Writing inode tables: done
Creating journal (32768 blocks): done
Writing superblocks and filesystem accounting information: done

[admin@localhost etc]$ sudo mount /dev/pmem0 /mnt/pmem/
[admin@localhost etc]$ df -h /mnt/pmem/
Filesystem      Size  Used Avail Use% Mounted on
/dev/pmem0      20G   45M   19G   1% /mnt/pmem
[admin@localhost etc]$
```

## 関連情報

- [クイック スタート ガイドインテル®オプタン™ DC永続メモリーのプロビジョニング](#)
- [永続メモリの設定](#)
- [インテル® Optane™ 永続メモリーの管理ユーティリティ—ipmctlおよびndctl](#)
- [テクニカル サポートとドキュメント – Cisco Systems](#)