Fujitsu Server PRIMERGY CX400 M4 Xeon スケーラブル・プロセッサ搭載システムの ための BIOS 最適化



本書では、Xeon スケーラブル・プロセッサ搭載 PRIMERGY CX400 M4 サーバ世代(PRIMERGY CX2550 M4 / M5、CX2560 M4 / M5、CX2570 M4 / M5)で設定可能な BIOS 設定について説明しています。

本書では、PRIMERGY サーバを使用するユーザが自身の要件に応じて BIOS 設定を最適化できるようにすることを目的としています。最適化の方向性としては、パフォーマンスとエネルギー効率のどちらかを最大化することを目指します。更にパフォーマンスについては、スループットを最大化するための最適化に加え、応答時間をできるだけ短くすることを重視するアプリケーションシナリオを検討します。

バージョン



目次

概要	3
アプリケーションシナリオ	
パフォーマンス	
低レイテンシ	
省エネ / エネルギー効率	
PRIMERGY BIOS オプション	
推奨される最適化設定	6
BIOS オプションの詳細	9
関連資料	20

概要

PRIMERGY サーバは、工場出荷時の時点で、最も一般的なアプリケーションシナリオ向けに、パフォーマンスとエネルギーの最もバランスのよい標準の BIOS 設定がなされています。ただし、可能な限り最大のスループット(パフォーマンス)、可能な限り最小のレイテンシ(低レイテンシ)、または可能な限り最大の省エネ(エネルギー効率)を重視したいという要件がある場合、標準設定からの変更が必要になる可能性があります。本書では、この3つのシナリオについて、最適な BIOS 設定として推奨されるベストプラクティスについて、以下に詳細に説明します。

PRIMERGY サーバを最適化する際は、BIOS 設定だけでなく、システム全体も考慮する必要があります。サーバシステムのプランニングには、次の点を特に考慮する必要があります。

・サーバハードウェア

・プロセッサ コア数および周波数

・メモリ メモリの種類(3DS DIMM、LR DIMM、RDIMM、NVDIMM)とメモリ構成

・I/O カード PCIe スロットにおける複数のカードの最適な配置

• オペレーティングシステムとアプリケーションソフトウェア

・ハイパーバイザ vSphere、Hyper-V、KVM

・プラン パフォーマンスやエネルギー効率

・チューニング カーネル、レジストリ、割り込みバインディング、スレッド分割

・ネットワーク

・ネットワークテクノロジ 1/10/25/40/100 Gbit イーサネット、ファイバーチャネル、InfiniBand、RDMA

・ネットワークアーキテクチャー スイッチ、マルチチャネル

・ストレージ

・テクノロジ RAID、ファイバーチャネル、Direct Attached、NVMe

・ディスク HDD、SSD、SATA、SAS

アプリケーションシナリオ



パフォーマンス

現在のオペレーティングシステムやアプリケーションに対応する最新のマルチプロセッサ、マルチコア、マルチスレッドテクノロジにより、Intel Xeon スケーラブル・プロセッサを搭載した今日の PRIMERGY サーバは、最高

レベルのパフォーマンスを提供します。これは、Standard Performance Evaluation Corporation (SPEC)、SAP 社、VMware 社による数々のベンチマーク性能の公開によっても証明されています。サーバのパフォーマンスにおいて重視されるのは、たいていはスループットについてです。最高のパフォーマンスを求めるユーザは、できるだけ多くの作業を同時に実行し、可能であれば並列プロセッサのすべてのリソースを活用したいと考えます。PRIMERGY サーバは、標準設定でもパフォーマンスとエネルギー効率に最もバランスのよい設定を提供しますが、BIOS 設定によって、システムのパフォーマンスとエネルギー効率を最大にするように最適化できます。パフォーマンスを最適化する場合はシステム内のすべてのコンポーネントを可能な限り最大速度で動作させ、省エネオプションの機能がシステムの速度低下を招かないようにします。そのため、パフォーマンスが最大になるように最適化すると、消費電力の増加につながります。



低レイテンシ

特にハイパフォーマンスコンピューティング(HPC)分野や、遅延なくリアルタイムで毎秒数百万のトランザクションとデータ処理を行う必要がある金融市場のアプリケーションでは、可能な限り最小のレイテンシが求められています。この分野のユーザは、システムの最適化を介して可能な限り最大のスループットを実現することでは

なく、個々のトランザクションの速度を上げること、すなわち、個々のトランザクションの実行にかかる時間を短縮することを 重視します。このような場合は、システムの応答時間、いわゆるレイテンシ(通常はナノ秒、マイクロ秒、またはミリ秒で測定)が 焦点になります。BIOS は、レイテンシを改善するさまざまなオプションを提供します。また、対応するアプリケーションがハー ドウェアで使用可能なすべてのスレッドを効率的に使用するわけではないことがわかっている場合、不要なスレッド(ハイパース レッディング)またはコアを BIOS 設定で無効にすることで、HPC アプリケーションで特に多く発生する演算速度の変動を最小限 にすることも可能です。さらに、不要なコアを無効にすることで、特定の動作条件下での残りのコアのターボモードのパフォー マンスを向上させることができます。一方、できるだけ一定のパフォーマンスを必要とするシナリオもあります。この場合は、 ターボモードなどで周波数変動が発生するような設定を避け、応答時間を一定に保つことが必要になります。現世代の Intel プロ セッサは、先行世代よりも明らかに優れたターボモードパフォーマンスを実現していますが、ターボモードの最大周波数は、特 定の動作条件下では保証されません。このような場合は、ターボモードを無効にすると、周波数変動を回避できます。省エネ機 能は、可能な限り周波数や電圧を低くし、特定の機能ブロックおよびコンポーネントを無効にしてエネルギーを節約することが 目的ですが、応答時間に悪影響を及ぼすこともあります。つまり、より強い省エネモードにすればするほど、パフォーマンスが 低下します。さらに、省エネモードで一旦低下したプロセッサのパフォーマンスを最大に戻すには、一定の時間を必要としま す。そのため、特にトランザクションが保留されてアイドル状態が続いていた後や、システムの負荷が不規則に変動している場 合、システムのレイテンシ増加につながります。本書では、低レイテンシを重視する分野のユーザを対象として、システムレイ テンシが最小限になるように省エネモードを構成する方法について説明します。しかし、サーバのレイテンシ、特にアイドル状 態のレイテンシを最適化すると必ず、電力消費量が実質的に多くなります。

「性能」および「低レイテンシ」に関する注意:

I/O システムの最大スループットまたは最小レイテンシは、I/O に強く依存するアプリケーションに大きな影響を与える場合があります。I/O システムのスループットまたはレイテンシの値は、プロセッサに対しては異なる意味を持ちます。例えば I/O スループットは、I/O システムによって一定時間内に転送されるデータの量を意味します。最大 I/O スループットまたは最小 I/O レイテンシを達成するために、BIOS のプロセッサ最適化機能を、最大のコンピュータ処理速度(「パフォーマンス」)または「低レイテンシ」に設定する必要はありません。ほとんどの場合では、最適に設定された I/O コンポーネントとともに BIOS の標準設定を使用するのが最も適しています。そうすることで、これらのコンポーネントに対して、ほぼ例外なく可能な限りの最適化がなされます。ただし、特定のまれなケース(要件が非常に高い SSD など)では、目標値が達成できない場合があります。その解決策として、BIOS オプションの[Uncore Frequency Scaling]を[Maximum]に設定するか、BIOS オプションの[Utilization Profile]を設定します(詳細については、それぞれのセクションを参照してください)。



省エネ / エネルギー効率

最大のスループットと最小のレイテンシのためのシナリオ、つまりパフォーマンスを重視するのではなく、エネルギー消費がより重要である環境もあります。この場合、次の2つの選択がありえます。

1 つは、可能な限り消費電力を低く抑えるように BIOS のオプションを選択することです。これは、電力予算に限りがあり、パフォーマンスよりもラックやサーバ当たりの電力消費量の削減を重視しているデータセンターオペレーターなどのユーザに適したオプションになります。この方向で最適化を行う場合は、速度、つまりサーバのパフォーマンスを低下させるための設定を行うことになります。

もう 1 つは、スループットと消費電力が最良の比率になるようにサーバを設定することです。これは、ワット当たりのパフォーマンスが測定されたサーバで最適なエネルギー効率を実現する唯一の方法です。こうした最適化は特に、パフォーマンスの最大化を重視せず、総所有コストの最適化に重きを置くデータセンターオペレーターによって採用されます。

Standard Performance Evaluation Corporation (SPEC)による数多くの公表と、サーバのエネルギー効率を測定する際の業界標準のベンチマークである SPECpower_ssj2008 や VMmark V3 Performance with Server Power は、PRIMERGY サーバがエネルギー効率に関しても最良の選択であることを証明しています。

PRIMERGY BIOS オプション

このホワイトペーパーには、Intel Xeon スケーラブル・プロセッサ搭載の PRIMERGY サーバに対して有効な BIOS オプションに関する情報が記載されています。これは以下のサーバに適用されます。

- PRIMERGY CX2550 M4 / M5
- PRIMERGY CX2560 M4 / M5
- PRIMERGY CX2570 M4 / M5

PRIMERGY サーバの BIOS は、常に開発が続けられています。そのため、いずれの場合も最新の BIOS バージョンを使用して、本書に記載されているすべての BIOS 機能を利用できるようにすることは大変重要です。該当するダウンロードは、https://www.fujitsu.com/global/support で公開されています。

推奨される最適化設定

以下の表に、最大のパフォーマンス、低いレイテンシ、または最大のエネルギー効率のいずれかを実現するために PRIMERGY サーバを最適化する場合の BIOS オプションの推奨設定を示します。BIOS オプションを変更するには、最初にシステムセルフテスト(Power On Self Test = POST)時の BIOS セットアップを呼び出す必要があります。詳細については、サーバのマニュアルを参照してください。

ここに記載されている BIOS オプションの多くは、互いに依存関係にあります。そのため、どのオプションの変更が、望ましくないシステムの動作を発生させ、また望ましいシステムの動作を発生させるかを明らかにするには、他のオプションも同時に変更してみるしかありません。以下の表に示されている BIOS オプションに変更する前に、該当の BIOS オプションの脚注に目を通すことをお勧めします。また、すべての変更を実稼働環境に適用する前に、必要な効果が有効かどうかテスト環境で検証することをお勧めします。

サーバシステムを計画する際は、BIOS オプションの推奨設定の他に、オペレーティングシステムの選択と調整にも留意する必要があります。使用方法によっては、オペレーティングシステムの選択と調整がパフォーマンス、レイテンシ、エネルギー効率に影響を及ぼす場合があります。個々のオペレーティングシステムの調整に関する補足情報については、「関連資料」の「オペレーティングシステムのパフォーマンス調整のガイドライン」のリンクを参照してください。

BIOS オプションの推奨最適化設定

ios セットアップメニュー	オプション設定1	パフォーマンス	低レイテンシ	エネルギー効率
onfiguration -> CPU Configuration				
Hyper-Threading	Disabled / Enabled	Enabled	Disabled ²	Enabled
Active Processor Cores	All / [1 - n]	All	1 - n ³	All
 Prefetcher Hardware Prefetcher Adjacent Cache Line Prefetch DCU Streamer Prefetcher DCU IP Prefetcher 	Disabled / Enabled	Enabled	Enabled	Disabled ⁴
 XPT Prefetch ⁵ LLC Prefetch ⁵ 	Disabled / Enabled	Enabled	Enabled	Disabled ⁴
Intel Virtualization Technology	Disabled / Enabled	Disabled ⁶	Disabled °	Disabled °
Power Technology	Disabled / Energy Efficient / Custom	Custom	Custom	Custom
Enhanced SpeedStep /	Disabled / Enabled	Enabled	Enabled	Enabled
Turbo Mode ^{7, 8}	Disabled / Enabled	Enabled	Disabled ⁹	Disabled
Energy Performance ^{7, 10}	Performance / Balanced Performance / Balanced Energy / Energy Efficient	Performance	Performance	Energy Efficien

¹ 太字で示している設定は標準値です。

² ハイパースレッディングを有効にすると論理コア数が 2 倍になりますが、性能のばらつきが発生することがあります。無効にすると、レイテンシが改善されます。

³ シングルスレッドのアプリケーション、またはすべての CPU スレッドを使用するわけではないアプリケーションでは、アクティブコアの数を制限するとターボモードのパフォーマンスが向上することがあります。

⁴ プリフェッチャーを無効にすると、パフォーマンスが変わらないか改善される場合にエネルギー効率が向上します。このことは、個々のプリフェッチャーについて前もって確認しておく必要があります。

⁵ M5 世代でのみ設定できます。

⁶ 仮想化を使用しない場合は、このオプションを[Disabled]に設定してください。

 $^{^{7}}$ [Power Technology]が Custom に設定されている場合のみ、このオプションが表示されます。

⁸ [Enhanced SpeedStep]が有効になっている場合のみ、このオプションが表示されます。

 $^{^9}$ すべての動作条件の下でターボモードの最大パフォーマンスが保証されているわけではないため、[Enabled]に設定するとパフォーマンスは変動します。[Turbo Mode]オプションを[Disabled]に設定すれば、安定した一定の応答時間にすることができます。

^{10 [}Override OS Energy Performance]の設定が[Enabled]に変更されている場合のみ、このオプションを設定できます。

BIOS セットアップメニュー	オプション設定 1	パフォーマンス	低レイテンシ	エネルギー効率
Override OS Energy Performance ^{7,}	Disabled / Enabled	Enabled	Enabled	Disabled ^{12^{エラー!} ブックマークが定義され ていません。}
Utilization Profile ^{7, 10}	Even / Unbalanced	Even	Unbalanced	Even
HWPM Support /	Disabled / Native Mode / OOB Mode / Native Mode with no legacy	Disabled	Disabled	Native Mode
CPU C1E Support /	Enabled / Disabled	Enabled	Disabled	Enabled
CPU C6 Report /	Disabled / Enabled	Enabled	Disabled	Enabled
Package C State limit [/]	C0 / C6 / No Limit	C0	C0	No Limit
UPI Link Frequency Select	Auto / 9.6 GT/s / 10.4 GT/s	Auto	Auto	9.6 GT/s
Uncore Frequency Scaling	Disabled / Enabled ¹³ / Nominal / Power balanced / Idle power saved	Disabled	Disabled	Nominal
Sub NUMA Clustering	Disabled / Enabled / Auto	Enabled	Enabled	Enabled
Stale AtoS	Disabled / Enabled	Enabled	Enabled	Enabled
LLC Dead Line Alloc	Disabled / Enabled	Disabled	Disabled	Disabled
Configuration -> Memory Configuration	on			1
DDR Performance	Performance optimized / Energy optimized	Performance optimized	Performance optimized	Energy optimized
Patrol Scrub	Disabled / Enabled	Enabled	Disabled	Enabled
DDR4 Write Data CRC Protection 5	Disabled / Enabled	Disabled	Disabled	Disabled

Fujitsu Public

^{11 [}HWPM Support]オプションが[OOB Mode]に設定されている場合このオプションは灰色表示され、その設定は自動的に[Enabled]に変更されています。

¹² 使用中のオペレーティングシステムで CPU の「エネルギー効率ポリシー」を設定することが可能な場合は、[Override OS Energy Performance]を[Disabled]に設定し、CPU の「エネルギー効率ポリシー」の設定をオペレーティングシステムの電源プランで行ってください。これが不可能な場合、またはオペレーティングシステムに設定させたくない場合は、このオプションを[Enabled]に設定し、BIOS の[Energy Performance]で設定を行ってください。

¹³ このオプションを[Enabled]に設定すると、I/O 稼働率は高いが、コア稼働率が低いアプリケーションで効果を発揮します。

BIOS オプションの詳細

本節では、各 BIOS オプションの詳細を記載します。

BIOS オプションを変更した場合の効果はハードウェア / ソフトウェア構成や、他の BIOS / OS オプション設定も影響するため、設定変更の際は、実運用の前に必ず検証を行ってください。

Hyper-Threading

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Hyper-Threading	Disabled Enabled	Enabled	Disabled	Enabled

通常、[Hyper-Threading]を有効([Enabled])にすることを推奨しています。ただし、金融市場の取引ソフトウェアや HPC アプリケーションのように、応答時間の短さを特に重要視するアプリケーションの場合は、[Hyper-Threading]を無効([Disabled])にすることをお勧めします。こうした分野のユーザは通常、追加スレッドによるシステムのスループットの最大化よりも、個々のスレッドのパフォーマンスと安定性を重視する傾向があります。[Hyper-Threading]を無効にするとパフォーマンス変動を抑えてレイテンシを改善することができます。

Active Processor Cores

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Active Processors Core	All [1 - n]	All	1 - n	All

BIOS でプロセッサの各コアを無効にすることもできます。たとえば、10 コアプロセッサ上で 4 つのコアを無効にできます。この場合、残りのコアは 10 コアの場合よりも多くの L3 キャッシュ容量を使用することができます。一般に最大のスループットは全てのコアを使用する場合に達成されますが、敢えて不要なコアを無効にすることで、残りのアクティブなコアでより高いターボモードの周波数の実現が期待できます。これは特に、すべてのコアを利用しない、レイテンシの影響を受けやすいアプリケーションの場合に有効です。不要なコアを無効にすることでプロセッサの電力消費量が低減され、残りのコアで実現可能なターボモード周波数が高くなります。これは、すべての負荷プロファイルで効果があるわけではありません。特に、消費電力の大きい AVX アプリケーションには効果が薄い場合があります。

Prefetcher

BIOS セットアップ メニュー	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
Configuration > CPU Configuration	Hardware Prefetcher Adjacent Cache Line Prefetch DCU Streamer Prefetcher DCU Ip Prefetcher	Disabled Enabled	Enabled	Enabled	Disabled
	XPT Prefetch LLC Prefetch	Disabled Enabled	Enabled	Enabled	Disabled

PRIMERGY サーバの BIOS には、上記の表に記載したプリフェッチャーオプションがあります。

プリフェッチャーはプロセッサの機能のひとつで、特定のパターンに応じてデータをメインメモリからプロセッサの L1 または L2 キャッシュに前もってロードすることができます。プリフェッチャーを有効にすると、通常、より高いキャッシュヒット率 を実現し、システム全体のパフォーマンスが向上します。これはメモリ転送がパフォーマンスのボトルネックになっているア プリケーションシナリオには向いていません。この場合、プリフェッチャーオプションを[Disabled]に設定して、プリフェッチ に使用される帯域幅をプリフェッチ以外に使用できるようにすることも可能です。また、プリフェッチャーを無効にすることで、サーバの消費電力をわずかに低減できます。実稼働システムでプリフェッチャーオプションを変更する前に、まずテスト 環境で各アプリケーションシナリオの個々の設定の効果を検証することをお勧めします。

個々のプリフェッチャーの詳細は以下の通りです。

Hardware Prefetcher	このプリフェッチャーは、データがアドレス A および A+1 で要求された場合、アドレス A+2 でも要求されることを想定してデータストリームを検索します。このデータはその後、メインメモリから L2 キャッシュにプリフェッチされます。
Adjacent Cache Line Prefetch	このプリフェッチャーは、データがキャッシュに格納されていなければ、常にキャッシュラインのペア(128 バイト)をメインメモリからキャッシュに転送します。このプリフェッチャーが無効の場合、プロセッサが要求するデータを含む 1 つのキャッシュライン(64 バイト)のみがキャッシュに転送されます。
DCU Streamer Prefetcher	このプリフェッチャーは、L1 データキャッシュのプリフェッチャーで、一定の時間内 に発生した同じキャッシュラインからの複数のロード命令を検出します。次のキャッシ ュラインも必要になるという仮定に基づいて、次のキャッシュラインが L2 キャッシュ またはメインメモリから L1 キャッシュに事前にプリフェッチされます。
DCU lp Prefetcher	この L1 キャッシュプリフェッチャーは、連続して実行されるメモリアクセス命令を検索し、次にアクセスされるデータを予測します。そして、必要に応じてこのデータを L2 キャッシュまたはメインメモリから L1 キャッシュにプリフェッチします。
XPT Prefetch	このプリフェッチャーは、LLC へのアクセスと並行してローカルメモリに対してもプリフェッチを行います。LLC でキャッシュミスした場合はプリフェッチしたデータを使用することでレイテンシを短縮します。プリフェッチ先の予測には、過去のアクセス履歴を基にした Xtended Prediction Table (XPT)を使用します。
LLC Prefetch	Xeon Scalable プロセッサでは、L3 キャッシュ(LLC: Last Level Cache)は non-inclusive キャッシュのため、通常はメインメモリから直接 L2 キャッシュにデータをリードしますが、このプリフェッチャーはメインメモリから L3 キャッシュにプリフェッチします。

Intel Virtualization Technology

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Intel Virtualization Technology	Disabled Enabled	Disabled	Disabled	Disabled

この BIOS オプションは、CPU の追加の仮想化機能を有効または無効にします。サーバを仮想化用に使用していない場合は、このオプションを[Disabled]に設定してください。これにより、電力を節約することもできます。

Power Technology

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Power Technology	Disabled Energy Efficient Custom	Custom	Custom	Custom

この BIOS オプションでは、CPU の電源管理機能の方針を設定します。[Disabled]に設定した場合、CPU の電源管理機能が無効になります。[Energy Efficient]に設定した場合、CPU の電源管理機能が省電力のために最適化されます。[Custom]に設定した場合、CPU の電源管理機能を手動で設定するためのいくつかのメニューが表示されるようになります。

Enhanced SpeedStep

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Enhanced SpeedStep	Disabled Enabled	Enabled	Enabled	Enabled

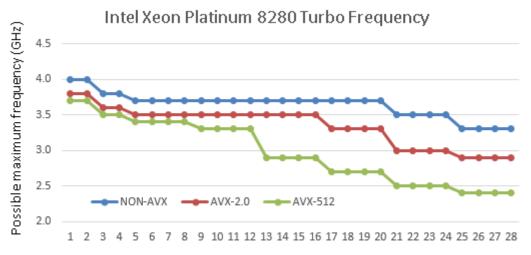
Enhanced Intel SpeedStep Technology (EIST)は、特定の負荷プロファイルに応じて各コアまたはプロセッサ全体のパフォーマンスを調整する省電力機能です。この機能では、最大演算速度が不要な場合に周波数と電圧を下げることで、必要なエネルギー量を大幅に下げます。演算速度の分散はオペレーティングシステムとオペレーティングシステムで導入されている戦略(実行する電力プラン)によって異なるため、富士通では[Enhanced SpeedStep]を有効にしたまま使用することをお勧めします。このオプションを無効にすると、ターボモード機能も利用できなくなります。

Turbo Mode

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Turbo Mode	Disabled Enabled	Enabled	Disabled	Disabled

この BIOS オプションは、プロセッサの Intel ターボブーストテクノロジ機能を有効または無効にします。ターボブーストテクノロジ機能では、周波数を公称周波数より上げることで演算速度を直ちに向上させることができます。実行可能な最大周波数は、プロセッサのタイプ、アクティブなプロセッサコアの数、電源、現在の電力消費量、温度、実行される命令(AVX512 命令か AVX2 命令か Non-AVX 命令か)など、多くの要素によって左右されます。

以下の図は、Xeon Platinum 8280 の実現可能な最大周波数を示しています。ここで、アクティブなプロセッサコア数とは、 [Active processor core]で有効にされ、かつ「C6 C-State」にないコアの数を意味します。(詳しくは[Active processor core]と [CPU C6 Report]の項を参照ください。)



Number of Active Processor Cores

これらの一般的な条件に加え、プロセッサの品質は、特に HPC アプリケーションの場合、ターボモードのパフォーマンスに大きく影響します。このため、たとえば同タイプのプロセッサ間の製造公差によっても、同負荷条件で電力消費量に違いが生じます。

通常は、[Turbo Mode]オプションを標準設定の[Enabled]に設定して、周波数を高くすることによりパフォーマンスを大きく向上させることを推奨しています。しかし周波数の高さとは上でも述べたように動作条件に依存し、常に保証されるものではないため、パフォーマンスを安定させたい、もしくは消費電力を少なくしたいアプリケーションシナリオでは、[Turbo Mode]を無効にすることをお勧めします。

Energy Performance

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Energy Performance	Performance Balanced Performance Balanced Energy Energy Efficient	Performance	Performance	Energy Efficient

この BIOS オプションは設定に応じて、Intel プロセッサ内部の「パワーコントロールユニット(PCU)」をパラメーター化して、プロセッサの電力管理機能をパフォーマンスとエネルギー効率の間で最適化します。可能な設定は、[Performance]、[Balanced Performance]、[Balanced Energy]、および[Energy Efficient]です。

Energy Performance 設定は、Energy Performance Bias とも呼ばれ、OS から設定することも可能ですが、BIOS オプション [Override OS Energy Performance]を[Enabled]に設定した場合には、BIOS で指定した本設定が強制的に有効になります。

[Override OS Energy Performance]を[Disabled]にした場合は、オペレーティングシステムが電源プランを介して[Energy Performance]オプションを設定するタスクを担いますが、OS の種類や設定によっては本設定が OS の電源ポリシーに影響する場合があります。

Override OS Energy Performance

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Override OS Energy Performance	Disabled Enabled	Enabled	Enabled	Disabled

新世代の Intel Xeon プロセッサは、さまざまな省エネオプションを備えています。プロセッサの中のいわゆるパワーコントロールユニット(PCU)は、これらの省エネオプションすべてを制御する際の中心的な役割をします。PCU は、設定を省エネや最大パフォーマンス重視で制御するために、パラメーター化できます。これには2つの方法があります。1つは[Override OS Energy Performance]の標準設定である[Disabled]設定を使用し、後述する[Energy Performance]オプションを、オペレーティングシステムを通じて制御する方法です。オペレーティングシステムで設定された電源プランに応じて、特定の値がCPU レジスターに書き込まれます。このレジスターを PCU が評価し、CPU の省エネ機能がそれに応じて制御されます。もう1つの方法は、[Override OS Energy Performance]を[Enabled]に設定することで、オペレーティングシステムの設定を無効にし、[Energy Performance] オプションを BIOS を介して直接設定します。これは特に、たとえば古いオペレーティングシステムでこの特殊な CPU レジスターに書き込めない場合、あるいは省エネオプションを BIOS で一元的に(つまりオペレーティングシステムとは無関係に)設定したい場合に有効です。

ハードウェア電源管理([HWPM Support])を[OOB Mode]で使用する場合は、[Override OS Energy Performance]オプションが標準で有効にされ、BIOS オプション[Energy Performance]で選択したエネルギー効率またはパフォーマンスに関する優先設定および PCU パラメーター設定が使用されます。

Utilization Profile

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Utilization Profile	Even Unbalanced	Even	Unbalanced	Even

[Utilization Profile]オプションは省エネオプションをパラメーター化するために使用します。このオプションは UPI と PCIe 帯域幅の両方をモニタリングして、使用率に基づいてプロセッサ周波数を適応させようとします。標準設定は[Even]ですが、これは CPU 負荷がすべてのプロセッサで均等に分散されていて、適切な周波数が CPU 使用率に基づいて最適に適合されていることが前提です。そのため[Even]設定では、プロセッサ周波数が積極的には増加しません。一方、[Unbalanced]設定は、CPU 負荷が低い場合に PCIe 使用率が高いアプリケーションシナリオを対象とします。GPGPU による構成がこの典型的な例です。その場合、オペレーティングシステムは CPU の使用率が低いことから低い周波数を要求しますが、実際には可能な最大 PCIe 帯域幅を実現するために高い周波数が必要になります。UPI または PCIe 使用率が高い場合、[Unbalanced]設定により、プロセッサの周波数は、CPU 使用率が低い場合でも、積極的に増大します。標準設定の[Even]の方がエネルギー効率が良いため、通常はこの設定にすることをお勧めします。しかし、高い PCIe 帯域幅が必要といったパフォーマンスの問題がアプリケーションシナリオに存在する場合は、[Unbalanced]設定によりこの問題が解消される可能性があります。

HWPM Support

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	HWPM	Disabled Native Mode OOB Mode Native Mode with no legacy	Disabled	Disabled	Native Mode

HWPM は HardWare Power Management(ハードウェア電力管理)の略で、Intel Broadwell プロセッサ世代に導入され、Skylake プロセッサ世代で機能が拡張された省電力機能です。HWPM を使用すると、Intel Enhanced SpeedStep テクノロジに基づいた従来の電源管理と同様の方法で、CPU 使用率に応じて、プロセッサ周波数の制御を想定する 2 つの動作モードを構成できます。従来の電源管理では、CPU 使用率の評価および P-State の制御がオペレーティングシステムによって、すなわちソフトウェアで制御されましたが、これとは対照的に、HWPM の場合は、これらのタスクはプロセッサ自体によってハードウェアで実施されます。HWPM は、従来の電源管理サポートを提供しないか、または従来の電源管理サポートを提供するが効率の悪いオペレーティングシステムにおいては、特に適切な選択肢です。

[Native Mode]を設定すると、オペレーティングシステムには、HWPM を使用した電源管理に関する制限および情報の伝達に使用されるインターフェースが提供されます。伝達された制限および情報は、HWPM による制御の際に考慮されます。一方、[OOB Mode]を設定すると、ハードウェア電源管理機構が自律的に、つまりオペレーティングシステムとは完全に独立して、プロセッサ周波数の制御を担います。[Native Mode with No Legacy]を設定すると、オペレーティングシステムに対して電源管理に関する制限および情報を HWPM の[Native Mode]で使用するインターフェースのみが提供されます。つまり、この設定では従来の P-state のインターフェースを提供しません。[HWPM Support]オプションが[Disabled]の場合には、HWPM を使用せず、[Enhanced SpeedStep]を使用した従来の電源管理が有効になります。

CPU C1E Support

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	CPU C1E Support	Enabled Disabled	Enabled	Disabled	Enabled

Xeon スケーラブル・プロセッサは、CO、C1、C1E、C6の4つのCステートをサポートします。C0以外のCステートはある種のスリープ状態です。消費電力は、C0、C1、C1E、C6の順で小さくなりますが、同じ順番でスリープ状態から復帰するのにかかる時間は長くなります。

C ステートの遷移は OS からのリクエストで発生します。このオプションを有効に設定した場合、C1 への遷移リクエストは CPU によって C1E への遷移リクエストとして扱われ、結果として消費電力がわずかですがさらに減少します。OS によって は、直接 C1E への遷移をリクエストすることもあり、この場合にはこのオプションは効果がありません。

C1E ステートでは、周波数のクロック数がそのプロセッサでサポートされる最低周波数に下がります。これは、Intel SpeedStep テクノロジに関係なく行われます。言い換えると、プロセッサが最大周波数で動作する設定がオペレーティングシステムの電源プランを介して行われていても、C1E が有効であれば、プロセッサはアイドル状態になるとクロック数が最低周波数に下がります。これは、特に低レイテンシアプリケーションで不利になる可能性があります。なぜなら、周波数のクロック数低下と復帰でレイテンシが増加するためです。そのような場合は、この設定を[Disabled]に変更できます。低レイテンシが必要とされるワークロードを除き、このオプションは[Enable]に設定することをお勧めします。

CPU C6 Report

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	CPU C6 Report	Disabled Enabled	Enabled	Disabled	Enabled

この BIOS オプションは、オペレーティングシステムに C6 ステートを使用できる([Enabled])か、あるいは使用できない ([Disabled])かを知らせるために使用されます。

この C ステートからの起床時間によってレイテンシが増加するため、可能な限り低いレスポンス時間での最大パフォーマンスが重要になるアプリケーションでは CPU C6 Report を[Disabled]にすることをお勧めします。C6 ステートが無効になると、最高のターボモード周波数を実現できなくなることを留意しておく必要があります。この場合はアクティブなコア数に関係なく、最高のターボモード周波数は、すべてのコアがアクティブな場合に可能となる最大周波数に限定されます。プロセッサタイプにもよりますが、これは通常はかなり小さくなります。ターボモード周波数が最大になるためには、すべてのコアが有効でない限り、[CPU C6 Report]を[Enabled]に設定する必要があります。BIOS オプション[CPU C6 Report]で[Disabled]設定を使用することによって BIOS ができるのは、ACPI を介して適切な CPU C ステートをオペレーティング システムへ転送することをやめることだけです。CPU コアの C ステートに関する BIOS 設定は、一部のオペレーティングシステム、特に「intel_idle」ドライバを使用する Linux ディストリビューションには効果がありません(2021 年現在、富士通がサポートしているすべての Enterprise Linux ディストリビューションが対象)。そのようなオペレーティングシステムで C ステートを設定する方法は 2 つあります。一つ目の方法は、BIOS で適切な C ステートの設定をして、Linux カーネルパラメーター「intel_idle.max_cstate=0」を使用してこのドライバを無効にすることです。こうすると、Linux カーネルは代わりに acpi 標準のアイドルドライバを使用するようになるため、BIOS 設定が有効になります。もう一つの方法は、Linux のコマンド「cpupower」を使用することです。このコマンドは BIOS 設定に関係なく期待する C ステートを設定することができます。

参考:プロセッサの電力状態



プロセッサパフォーマンス 電力状態(P-State)

- Enhanced Intel SpeedStep Technology (EIST)や Demand Based Switching (DBS)と呼ばれる
- P-State は、プロセッサがコードを実行している場合 にも、CPU 使用率に基づいて消費電力を小さくする
- P-State はプロセッサ電圧とプロセッサ周波数で制御 される
- P-State は、さまざまなパフォーマンスレベルがある



プロセッサアイドル時 電力状態 (C-State)

- C-State は、プロセッサがコードを実行していない場合に消費電力を少なくする
- プロセッサの一部を無効にできる
- CO → プロセッサがアクティブ
- C6 → プロセッサがディープパワーダウン状態

Package C State limit

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Package C State limit	C0 C6 No Limit	C0	C0	No Limit

CPU またはコア C ステートに加えて、いわゆるパッケージ C ステートというものもあります。このときは、個々のコアのプロセッサのみでなく、プロセッサチップ全体をある種のスリープ状態にすることができます。その結果として、消費電力量はさらに少なくなります。低パッケージ C ステートからアクティブな CO ステートへ変わるのに必要な「ウェイクアップ時間」は、CPU またはコア C ステートと比べると長くなります。[C0]設定が BIOS で行われると、プロセッサチップは常にアクティブなままになります。ただし、動作時間内のサーバのアイドル時間が非常に長くなることが予想できず、パッケージ C ステートからの「ウェイクアップ」時にレイテンシが重要な役割を果たさない場合は、アイドル状態のサーバの電力消費量を大幅に低減するために、設定を[C6]にします。

UPI Link Frequency Select

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	UPI Link Frequency Select	Auto 9.6 GT/s 10.4 GT/s	Auto	Auto	9.6 GT/s

この BIOS オプションを使用することで、システム内のプロセッサ間をつなぐ UltraPath インターコネクト(UPI)の動作周波数を小さくして電力を節約することができます。これは、UPI の使用可能帯域幅が不要な場合に特に有効です。しかしパフォーマンスを最大にしてレスポンス時間を短くすることが指定されている場合は、最高速度を自動的に設定する[Auto]設定のまま変更しないでおきます。必要となる帯域幅に応じて、最大限の省電力を実現する[9.6 GT/s]、最高速度である[10.4 GT/s]から選択できます。

Uncore Frequency Scaling

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Uncore Frequency Scaling	Disabled Enabled Nominal Power balanced Idle power saved	Disabled	Disabled	Nominal

Intel Xeon スケーラブル・プロセッサは、個々のコア、そしてアンコアと呼ばれる領域が、それぞれ独立した周波数で動作します。各領域の周波数は、稼働率に応じて設定されます。つまり、稼働率が高くなるとプロセッサの周波数が高くなり、負荷に応じた演算性能を発揮できます。一方で、プロセッサやプロセッサの該当する領域で稼働率が低下すると、周波数を最小限に抑えてエネルギーを節約します。

この BIOS オプションの設定は、アンコア領域の周波数を制御します。[Disabled]を使用すると、アンコア周波数がプロセッサ自体によって制御されます。アンコア周波数は、現在の CPU 使用率に従って、最低アンコア周波数と、最大アンコア周波数の間で変動します。アンコア周波数は、使用しているプロセッサの種類によって異なり、結果的にプロセッサの公称周波数を上回ることも下回ることもあります。標準設定の[Enabled]設定にすると、コアの稼働率が低い場合やアイドル状態の場合でも、プロセッサのアンコア領域が常に最大周波数で動作するようになりますが、それに応じて、電力消費量も高くなります。この

ため、通常はこのオプションを常に[Disabled]に設定します。I/O レイテンシが重要なアプリケーションや、一般にI/O を多用するアプリケーションは、プロセッサに対する負荷がまったくないか非常に少ないので、例外となります。この場合、プロセッサの電源管理メカニズムが周波数を最小に設定しようと試みます。この場合は、アンコア領域の周波数も自動的に低くなります。これは、I/O スループットに悪影響を与える可能性があります。というのも、すべてのI/O 通信(PCIe、メモリ、UPI など)はアンコア領域を経由しているためです。[Uncore Frequency Scaling]を[Enabled]に設定すればこの動作を回避できますが、電力消費量の増加は避けられません。[Nominal]設定は、最大可能アンコア周波数が最大でもプロセッサの公称周波数に制限されることを除いては、標準設定の[Disabled]と同じように動作します。I/O 稼働率が低いアプリケーションでは、これにより、エネルギー効率が向上することがあります。[Power balanced]設定では、電力消費量と性能が最適なバランスになるように動作します。[Idle power saved]設定では、システムがアイドル状態のときにアイドル電力を節約するようにプロセッサがアンコア周波数を動作します。

Sub NUMA Clustering

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Sub NUMA Clustering	Disabled Enabled Auto	Enabled	Enabled	Enabled

Sub NUMA Clustering (SNC)の設定は、L3 キャッシュをアドレス範囲に応じて 2 つのクラスタに分割します。

いずれのクラスタも、どちらか一つのメモリコントローラにくくりつけられます。また、オペレーティングシステムからは一つの NUMA ドメインとして扱われ、NUMA ノード内の L3 キャッシュやメモリのアクセスは、そのレイテンシが改善します。

SNC は、前の世代の CPU であった Cluster on Die (COD)の代替です。COD と同じように、SNC はローカルメモリレイテンシを最小化、ローカルメモリ帯域を最大化することができるため、NUMA 最適化されたアプリケーションにおいて特に推奨されます。

Stale AtoS

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	Stale AtoS	Disabled Enabled	Enabled	Enabled	Enabled

Xeon スケーラブル・プロセッサでは、I、A、S の 3 つのメモリディレクトリ状態があります。I (Invalid) 状態は、データがクリーン(メモリからデータを読み込んだ後にキャッシュ上で更新されていない)で、どの CPU のキャッシュにもそのデータが含まれていないことを意味します。A (SnoopALL)状態は、データが Exclusive または Modified の状態で存在する可能性があることを意味します。S (Shared)状態は、データがクリーンで、複数の CPU のキャッシュで共有されている可能性があることを意味します。

メモリリードを行う際、ディレクトリが A 状態の場合には、すべての CPU をスヌープしなければなりません。なぜなら、他の CPU がデータをキャッシュ上で更新している可能性があるためです。この場合、スヌープを行うと更新されたデータがキャッシュから転送されてきます。しかし、ディレクトリが A 状態と読めた場合でも、スヌープがミスとなることがあります。これ は、他の CPU がデータをキャッシュ上に読み込んだ後に、そのデータを更新することなく無言で廃棄した場合に起こります。 [Stale AtoS]が[Enabled]の場合には、A 状態のキャッシュラインへのスヌープがミスになったとき、そのラインの状態が S 状態に遷移します。こうすることで、それ以降のリードが発生した場合、S 状態であるためスヌープが不要となり、レイテンシとスヌープの帯域を削減することができます。[Stale AtoS]は、リモート CPU からのリードが多いワークロードで有効な可能性があります。

LLC Dead Line Alloc

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > CPU Configuration	LLC Dead Line Alloc	Disabled Enabled	Disabled	Disabled	Disabled

Xeon スケーラブル・プロセッサのキャッシュ機構では、L2 キャッシュから追い出されたデータが、L3 キャッシュに格納されます。L2 キャッシュからキャッシュラインが追い出されると、コアはそのキャッシュラインに「dead」というフラグを付けることがあります。「dead」とは、再びリードされる可能性が低い、ということを意味します。

[LLC Dead Line Alloc]が[Disabled]の場合には、「dead」というフラグのついたキャッシュラインが L3 キャッシュに格納されることはありません。これにより L3 キャッシュの領域が節約され、必要とするデータが L3 キャッシュから追い出されることを避けることができます。[LLC Dead Line Alloc]が[Enabled]の場合には、L3 キャッシュに空きがあれば、「dead」といフラグのついたキャッシュラインを L3 キャッシュに格納することがあります。

比較測定の結果、整数演算のワークロードにおいて[LLC Dead Line Alloc]が[Disabled]の場合に、わずかながら性能が高いことがわかりました。ただし、この効果はアプリケーションのキャッシュの使い方に大きく依存するため、このオプションを変更するまえにその効果をテスト環境で検証することをお勧めします。

DDR Performance

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > Memory Configuration	DDR Performance	Performance optimized Energy optimized	Performance optimized	Performance optimized	Energy optimized

この BIOS オプションは、メモリモジュールが動作する速度を制御します。このため、パフォーマンスとエネルギー消費量を比較しながら評価する必要があります。[Performance optimized]設定では、使用する CPU のタイプとメモリ構成に応じて DIMM が最大速度で動作するため、最高のメモリパフォーマンスが得られますが、消費電力も上がるため、性能と電力のトレードオフが発生します。[Energy optimized]設定では、プロセッサのモデルやメモリ構成にかかわらず、メモリ周波数が搭載したプロセッサでサポートされる最小周波数に常に制限されるため、電力消費量が削減されます。

Patrol Scrub

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > Memory Configuration	Patrol Scrub	Disabled Enabled	Enabled	Disabled	Enabled

この BIOS オプションは、システムのメインメモリに、オペレーティングシステムとは無関係にバックグラウンドで繰り返しアクセスして、メモリエラーを予防的に検出して修正する、いわゆるメモリスクラビングを有効または無効にします。一般的なワークロードでは[Patrol Scrub]オプションを有効にしても性能影響は小さいものの、このメモリテストの時間は調整することができず、特定の状況ではパフォーマンスがばらつく原因になる可能性があります。[Patrol Scrub]オプションを無効にすることにより、オペレーティングシステムによるアクセスがアクティブな場合にメモリエラーが検出される可能性が高まります。これらのエラーが修正可能な頻度で起きている限りは、メモリモジュールの ECC テクノロジによりシステムは引き続き安定して動作します。しかし、修正可能なメモリエラーが多すぎると、修正不可能なエラーが検出されるリスクが高まり、その結果としてシステムが停止してしまいます。

DDR4 Write Data CRC Protection

BIOS セットアップ	BIOS オプション	設定	パフォーマンス	低レイテンシ	エネルギー効率
メニュー					
Configuration > Memory Configuration	DDR4 Write Data CRC Protection	Disabled Enabled	Disabled	Disabled	Disabled

この BIOS オプションでは、DDR4 の Write CRC 機能を制御します。この設定を有効にすると、CPU 内蔵のメモリコントローラは DRAM への書き込み時に、書き込みデータとともに生成した CRC 符号を DRAM に送信します。DIMM 側で CRC をチェックし、データバスの 1bit エラー、2bit エラー、奇数 bit や水平カラムのマルチビットエラーを検出することができます。メモリパスの信頼性が向上する一方で、CRC を生成するためにレイテンシが悪化し、かつ余分にデータバスを使用するため、メモリ帯域が悪化します。

関連資料

PRIMERGY サーバ

https://www.fujitsu.com/jp/products/computing/servers/primergy/

PRIMERGY CX400 M4 Xeon スケーラブル・プロセッサ搭載システムのための BIOS 最適化

このホワイトペーパー

- http://docs.ts.fujitsu.com/dl.aspx?id=e693e228-33b5-48f7-8310-21a313191103
- http://docs.ts.fujitsu.com/dl.aspx?id=88f451a3-20ee-4989-a0ec-e21ec824ecd3

PRIMERGY のパフォーマンス

https://jp.fujitsu.com/platform/server/primergy/performance/

PRIMERGY のマニュアル

https://www.fujitsu.com/jp/products/computing/servers/primergy/manual/

サポートページ:

https://support.ts.fujitsu.com/

"BIOS Setup Utility"は機種ごとの以下のドキュメント名を検索することでダウンロードできます。

・ CX400 M4 / M5: "D385x BIOS セットアップユーティリ ティ"

オペレーティングシステムのパフォーマンス調整のガイドライン

· Microsoft Windows:

https://docs.microsoft.com/en-us/windows-server/administration/performance-tuning/

• RedHat Linux:

https://access.redhat.com/documentation/en-us/red hat enterprise linux/7/html/performance tuning guide/index https://access.redhat.com/documentation/en-

<u>us/red_hat_enterprise_linux/8/html/monitoring_and_managing_system_status_and_performance/index_https://access.redhat.com/documentation/en-</u>

us/red hat enterprise linux/9/html/monitoring and managing system status and performance/index

SUSE Linux:

https://documentation.suse.com/sbp/all/html/SBP-performance-tuning/index.html

https://documentation.suse.com/sles/15-SP3/html/SLES-all/book-tuning.html

https://documentation.suse.com/sles/15-SP4/html/SLES-all/book-tuning.html

https://documentation.suse.com/sles/15-SP5/html/SLES-all/book-tuning.html

• VMware vSphere:

https://www.vmware.com/files/pdf/techpaper/VMW-Tuning-Latency-Sensitive-Workloads.pdf

https://www.vmware.com/techpapers/2019/vsphere-esxi-vcenter-server-67U2-performance-best-practices.html

https://www.vmware.com/techpapers/2022/tagging-vsphere70u1-perf.html

https://www.vmware.com/techpapers/2021/vsphere-esxi-vcenter-server-70U2-performance-best-practices.html https://www.vmware.com/techpapers/2022/vsphere-esxi-vcenter-server-70U3-performance-best-practices.html

https://www.vmware.com/techpapers/2022/vsphere-esxi-vcenter-server-80-performance-best-practices.html

 $\underline{https://www.vmware.com/techpapers/2023/vsphere-esxi-vcenter-server-80U1-performance-best-practices.html}$

文書変更履歴

版数	日付	説明	
1.2	2023-10-03	新 Visual Identify フォーマットに変更 軽微な修正	
1.1	2019-06-06	M5 世代向けの設定を追加	
1.0	2018-04-05	初版	

お問い合わせ先

富士通株式会社

Web サイト: https://global.fujitsu/ja-jp/

PRIMERGY のパフォーマンスとベンチマーク

mailto:fj-benchmark@dl.jp.fujitsu.com