



Red Hat Enterprise Linux 7

電力管理ガイド

Red Hat Enterprise Linux 7 での電力消費量の管理

Red Hat Inc.
Rüdiger Landmann

Jacquelynn East
Jack Reed

Don Domingo

Red Hat Enterprise Linux 7 での電力消費量の管理

Jacquelynn East
Red Hat Engineering Content Services

Don Domingo
Red Hat Engineering Content Services

Rüdiger Landmann
Red Hat Engineering Content Services

Jack Reed
Red Hat Engineering Content Services

Red Hat Inc.

法律上の通知

Copyright © 2013 Red Hat Inc..

This document is licensed by Red Hat under the [Creative Commons Attribution-ShareAlike 3.0 Unported License](#). If you distribute this document, or a modified version of it, you must provide attribution to Red Hat, Inc. and provide a link to the original. If the document is modified, all Red Hat trademarks must be removed.

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, JBoss, MetaMatrix, 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 Software Collections is not formally related to or endorsed by the official Joyent Node.js open source or commercial project.

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

All other trademarks are the property of their respective owners.

概要

本ガイドでは、Red Hat Enterprise Linux 7 システムで効果的に電力消費量を管理する方法を説明しています。以下のセクションでは、電力消費量を低減する様々な技術 (サーバー向けとノート PC 向けの両方)、そしてその各技術がどのようにシステムの全体的なパフォーマンスに影響を与えるかについて説明しています。

目次

第1章 概要	2
1.1. 電力管理の重要性	2
1.2. 電力管理の基礎	3
第2章 電力管理の監査と分析	5
2.1. 監査および分析の概要	5
2.2. PowerTOP	5
2.3. Diskdevstat と netdevstat	7
2.4. Battery Life Tool Kit	11
2.5. Tuned	12
2.6. UPower	22
2.7. GNOME の電源管理	23
2.8. 他の監査ツール	23
第3章 中核となるインフラストラクチャとメカニズム	24
3.1. CPU のアイドル状態	24
3.2. CPUfreq ガバナーの使用	24
3.3. CPU モニター	27
3.4. CPU 節電ポリシー	28
3.5. サスペンドと復帰	28
3.6. Active-State Power Management	28
3.7. Aggressive Link Power Management	29
3.8. Relatime ドライブアクセス最適化	30
3.9. パワーキャッピング (Power Capping)	30
3.10. 拡張グラフィックス電力管理	31
3.11. RFKill	32
第4章 使用例	34
4.1. 例 — サーバー	34
4.2. 例 — ノート PC	35
付録A 開発者へのヒント	37
A.1. スレッドの使用	37
A.2. ウェイクアップ (Wake-ups)	38
A.3. Fsync	39
付録B 改訂履歴	41

第1章 概要

Red Hat Enterprise Linux 7 の重要な改善点の 1 つが電力管理です。コンピューターシステムで使用する電力を制限することは、グリーン IT (環境に優しいコンピューティング) の最も重要な側面の 1 つです。このグリーン IT では、リサイクル可能な資源の利用、ハードウェアの製造により環境に及ぼす影響、システム設計と導入における環境意識などについても考慮されます。本ガイドでは、Red Hat Enterprise Linux 7 が稼働するシステムの電力管理について説明します。

1.1. 電力管理の重要性

電力管理の中核は、各システムコンポーネントによるエネルギーの消費をいかに効果的に最適化するかを理解するところにあります。そのためには、システムによって行なわれるさまざまなタスクを調査し、そのパフォーマンスがジョブにぴったりと適したパフォーマンスとなるよう各コンポーネントの設定を行なっていく必要があります。

電力管理を行う主な要因を以下に示します。

- ▶ コスト削減のため全体的な消費電力を節電する

電力管理を適切に活用すると、以下のような結果が得られます。

- ▶ サーバーおよびコンピューティングセンターの冷却
- ▶ 冷却、空間、ケーブル、発電機、無停電電源装置 (UPS) などにかかる二次コストの削減
- ▶ ノートパソコンのバッテリー寿命の延長
- ▶ 二酸化炭素排出量の低減
- ▶ エナジースター (Energy Star) などグリーン IT に関する政府の規則、又は法的要件への適合
- ▶ 新システムにおける企業のガイドラインへの適合

通例、あるコンポーネント (またはシステム全体) の電力消費を抑えようとする、発生熱量が低下するため、必然的にパフォーマンスの低下につながります。そのため、特にミッションクリティカルなシステムの場合、設定を変更することで得るパフォーマンスの低下については十分な調査と検証を行なってください。

システムにより行なわれる様々なタスクを調査し、そのパフォーマンスがジョブにぴったりと適したパフォーマンスとなるよう各コンポーネントの設定を行なうことで、エネルギーを節約し、発生熱を低減、ノートパソコンのバッテリー寿命を最適化することができます。電力消費に関するシステムの分析とチューニングに関する原則の多くは、パフォーマンスのチューニングの原則と似ています。通常、システムはパフォーマンスまたは電力のいずれかに対して最適化が行なわれるため、電力管理とパフォーマンスのチューニングは、ある意味、システムを構成する上で正反対となるアプローチになると言えます。本ガイドでは、電源管理を行なう上で役に立つ Red Hat 提供のツール、そして Red Hat で開発された技術について説明していきます。

Red Hat Enterprise Linux 7 には、デフォルトで有効になっている多くの新しい電力管理機能がすでに含まれています。これらはすべて、サーバーやデスクトップの標準的な使用でパフォーマンスに影響を及ぼさないように選択されています。ただし、最大のスループット、最小限の遅延、または最大の CPU パフォーマンスが必要な非常に特殊なユースケースでは、これらのデフォルト値を再検討する必要がある場合があります。

以下の質問とその答えをお読みいただいた上で、本ガイドで説明している技術を使用してマシンを最適化すべきかどうかの判断を行なってください。

問：

マシンを最適化すべきですか？

答：電力を最適化する重要性は、お客様に従うべきガイドラインがあるか、又は順守すべき規則があるかによって変わってきます。

問：

どの程度、最適化する必要がありますか？

答：ここで紹介している技術の中には、マシンを詳しく監査、分析する行程すべてを通して行なう必要がなく、その代わりに電力の使用を標準的に改善する全般的な最適化を行えばよいものがあります。もちろん手作業で行うシステム監査と最適化ほど優れてはいませんが、適度な効果はあります。

問：

最適化によりシステムパフォーマンスが許容範囲を下回るレベルまで低減されてしまいましたか？

答：本ガイドで説明しているほとんどの技法により、システムパフォーマンスは明らかな影響を受けません。Red Hat Enterprise Linux 7 にすでに設定されているデフォルト値を使用せずに電力管理を実装する場合は、電力最適化の実行後にシステムパフォーマンスを監視し、パフォーマンスの低下が許容範囲内であるか判断する必要があります。

問：

最適化に時間とリソースを費した場合、そこから得られる結果より負担の方が大きくなってしまいましたか？

答：1台のシステムに対して全行程を手作業で行なっていく最適化については、費される時間とコストが1台のマシンの寿命が尽きるまでに得られるであろう恩恵をはるかに上回ってしまうため、一般的には意味がありません。一方、例えば1万台のデスクトップシステムに同じ構成と設定を持たせて、複数のオフィスへの実装を展開する場合には、最適化した設定をひとつ構成してそれを1万台すべてのマシンに適用していけば十分に役立つ可能性が高くなります。

次のセクションでは、最適なハードウェアのパフォーマンスがエネルギー消費の観点でどのようにシステムに恩恵をもたらすのかについて解説していきます。

1.2. 電力管理の基礎

効率的な電力管理は以下の原則の上に成り立っています。

アイドル状態の CPU は必要な時にウェイクアップする

Red Hat Enterprise Linux 6 以降、カーネルは、ティックレス (*tickless*) を実行します。つまり、旧式の定期タイマーの割り込みが、オンデマンド型の割り込みに取って代わられます。そのため、アイドル状態の CPU を、新しいタスクが処理のためにキューに格納されるまでアイドル状態に維持し、低電力の状態にある CPU をより長くその状態に維持することができます。ただし、システムのアプリケーションが不必要なタイマーイベントを作成する場合、この機能の利点は打ち消されることがあります。

Red Hat Enterprise Linux 7 には、CPU の使用量でアプリケーションを識別し、監査を行なうことができるツールが同梱されています。詳細は、[2章 電力管理の監査と分析](#) を参照してください。

使用していないハードウェアとデバイスは完全に無効にする

これは、可動パーツを持つデバイス (例えば、ハードディスク) に特に当てはまります。さらに、一部のアプリケーションは、使用していないが有効なデバイスを「オープン」状態にすることがあります。これが起こると、カーネルはデバイスが使用中だと想定し、デバイスが節電状態に入ることを阻止する可能性があります。

動作が少ないということは消費電力も少ない

ただし多くの場合、これは最新のハードウェアと正しい BIOS 設定によって異なります。現在 Red Hat Enterprise Linux 7 では対応できるようになった新しい機能の一部は、旧式のシステムコンポーネントでは対応していないことが多々あります。システムに最新の公式のファームウェアが使用されていること、また BIOS の電力管理またはデバイス設定のセクションで電力管理の機能が有効になっていること、を確認してください。確認する機能は以下のとおりです。

- ✦ SpeedStep
- ✦ PowerNow!
- ✦ Cool'n'Quiet
- ✦ ACPI (C 状態)
- ✦ Smart

上記の機能がハードウェアでサポートされ、BIOS で有効になっている場合、Red Hat Enterprise Linux 7 ではその機能がデフォルトで使用されます。

CPU の各種状態とその効果

ACPI (電力制御インタフェース : *Advanced Configuration and Power Interface*) を搭載する最新の CPU は、以下の 3 種類の電力状態を提供します。

- ✦ Sleep (C 状態)
- ✦ Frequency (P 状態)
- ✦ Heat output (T 状態、または「温度状態」)

最小限のスリープ (sleep) 状態で稼働している CPU は、最小限のワット数を消費しますが、必要なときにこの状態から復帰するには相当な時間がかかります。非常に稀なケースですが、CPU がスリープに入る度に直ちに復帰する必要がある場合もあります。この場合、事実上 CPU が常にビジーな状況を引き起こし、別の状態を使用していたら実現できた節電効果が得られなくなってしまうこととなります。

電源がオフになっているマシンの消費電力は最小となる

当たり前かもしれませんが、実際に節電を行う最善策の 1 つは、システムの電源を切ることです。例えば、「グリーン IT」の意識に焦点を置いた企業文化を育み、昼休みや帰宅時にはマシンの電源を切るガイドラインを設けるのも一案です。また、数台の物理サーバーを大きなサーバー 1 台に統合し、Red Hat Enterprise Linux 7 で配布されている仮想化技術を使用して仮想化することもできます。

第2章 電力管理の監査と分析

2.1. 監査および分析の概要

たった一台のシステムに対して監査や分析、チューニングを細かく手作業で行うことは、通常、例外的です。こうしたシステムのチューニングの作業にかかる時間やコストが、一般的にはその作業から得られる恩恵を上回ってしまうためです。しかし、この作業を一度だけ行い、同じ設定をほぼ同一構成となる大量のマシンに再利用できる場合には、非常に便利です。たとえば、数千に及ぶデスクトップシステムの導入、あるいはほぼ同一構成の複数マシンから成る HPC クラスターの導入などを考えてみてください。監査や分析を行うもうひとつの理由は、将来的にシステムの動作に起こる後退や変化を特定できるよう比較対象となる基準を設けるといことです。ハードウェアや BIOS、ソフトウェアなどの定期更新によって予想以上の電力消費が発生するのを避けたい場合などに、この分析結果が非常に役立ちます。一般的には、徹底的な監査や分析を行うことで、特定のシステムで発生している現状を把握できるようになります。

電力消費に関する監査と分析は、最新システムを使用しても比較的難しいものです。ほとんどのシステムは、ソフトウェアを介する電力使用量を測定するために必要な手段を提供しません。ただし、例外はあります。Hewlett Packard サーバーシステムの ILO 管理コンソールには、ウェブ経由でアクセスできる電力管理モジュールが備わっています。IBM は、BladeCenter 電力管理モジュールで同様のソリューションを提供しています。Dell システムの一部でも、IT Assistant 機能により電力監視機能が提供されています。他のベンダーはサーバープラットフォーム向けに類似の機能を提供している可能性がありますが、すべてのベンダーで対応しているソリューションは存在しません。

電力消費を測定する直接的な方法としては多くの場合、できるだけ節電に最大限の努力をすることしかありません。変更が反映されているか、システムがどのように動作しているか査定する他の方法もあります。この章では、そのために必要なツールについて説明します。

2.2. PowerTOP

ティックレスカーネルでは、CPU を頻繁にアイドル状態に切り替え、電力消費を削減して電力管理を改善することができます。新しい **PowerTOP** ツールを使用すると、CPU を頻繁にウェイクアップするカーネルの特定のコンポーネントとユーザースペースアプリケーションを特定できます。

Red Hat Enterprise Linux 7 には、バージョン 2.x の **PowerTOP** が含まれます。このバージョンは、1.x コードベースの完全な書き換えであり、わかりやすいタブベースのユーザーインターフェースを備え、カーネルの "perf" インフラストラクチャーを広範に使用してより正確なデータを提供します。システムデバイスの電力動作が追跡され、明確に表示されるため、問題を迅速に特定することが可能です。試験的に、2.x コードベースには、個別のデバイスおよびプロセスが消費している電力を示すことができる電力予測エンジンが含まれています。[図2.1「実行中の PowerTOP」](#)を参照してください。

PowerTOP をインストールするには、**root** で以下のコマンドを実行します。

```
yum install powertop
```

PowerTOP を実行するには、**root** で以下のコマンドを実行します。

```
powertop
```

PowerTOP はシステム全体の電力使用量の予測を提供し、各プロセス、デバイス、カーネル作業、タイマー、割り込みハンドラーの電力使用量を表示することができます。このタスク中は、ノート PC をバッテリー電源で稼働してください。電力予測エンジンを調整するには、**root** で以下のコマンドを実行します。

```
powertop --calibrate
```

調整には時間がかかります。このプロセスでは様々なテストが実行され、輝度レベルおよびスイッチデバイスのオンとオフが繰り返されます。調整中にはマシンに触れないでください。調整プロセスが終わると、**PowerTOP** が正常に開始されます。データ収集をおよそ 1 時間実行させます。十分なデータが収集されると、最初のコラムに電力予測の数字が表示されます。

ノート PC でこのコマンドを実行する場合は、バッテリー電源で稼働することで利用可能なデータすべてが提供されます。

PowerTOP は実行中にシステムから統計数字を収集します。**Overview** タブでは、CPU にウェイクアップを最も頻繁に送信するコンポーネントまたは最も電力を消費しているコンポーネントのリストが表示されます (図 2.1 「実行中の PowerTOP」を参照)。その横のコラムでは、電力消費予測、リソースの使用方法、1 秒あたりのウェイクアップ、プロセスやデバイス、タイマーなどコンポーネントの分類、およびコンポーネントの説明が表示されます。1 秒あたりのウェイクアップは、サービスまたはカーネルのデバイスおよびドライバーのパフォーマンスの効率性を示します。ウェイクアップが少ないと消費電力も少ないこととなります。コンポーネントは、電力使用量の最適化がさらに実行可能な度合いで並んでいます。

ドライバーコンポーネントのチューニングは通常、カーネルの変更を必要とし、本ガイドの対象外となります。ただし、ウェイクアップを送信するユーザー空間のプロセスは、管理がより簡単です。最初に該当するサービスまたはアプリケーションをこのシステム上で実行する必要があるかどうかを判断します。必要ない場合は、そのサービスまたはアプリケーションを単に非アクティブ化します。古い System V サービスを永続的に無効にするには、以下のコマンドを実行します。

```
systemctl disable servicename.service
```

このプロセスについてより詳細な情報を得るには、**root** で以下のコマンドを実行します。

```
ps -awux | grep processname  
strace -p processid
```

トレースが繰り返し行われているように見える場合は、恐らくビジーループが発生しています。このようなバグを修復するには通常、そのコンポーネントでコードを変更する必要があります。

図 2.1 「実行中の PowerTOP」では、消費電力量の合計と、該当する場合はバッテリーの残量が表示されます。これらの中には、1 秒あたりのウェイクアップ合計、1 秒あたりの GPU 操作、および 1 秒あたりの仮想ファイルシステム操作の概要があります。画面の残りには、使用量にしたがってプロセス、割り込み、デバイス、およびリソースの一覧が表示されます。適切に調整されると、一覧の各アイテムの最初のコラムに電力消費量の予測も表示されます。

タブを移動するには、**Tab** および **Shift+Tab** キーを使用します。**Idle stats** タブでは、すべてのプロセッサおよびコアの C 状態の使用が表示されます。**Frequency stats** タブでは、Turbo モード (該当する場合) を含む全プロセッサおよびコアの P 状態の使用が表示されます。CPU がより高い C または P 状態に長くいればいるほど、よいこととなります (C4 の方が C3 よりも高い)。これは、CPU 使用率がどの程度うまく最適化されているかを示す指標となります。システムのアイドル中の理想状態は、最高の C または P 状態が 90% 以上を維持していることです。

Device Stats タブは **Overview** タブと同様の情報を表示しますが、デバイスに限定されます。

Tunables タブには、システムの消費電力量を低減させるための提案が含まれています。**up** および **down** キーを使って各提案に移動し、**enter** キーでそれらのオン/オフを切り替えます。

```

PowerTOP 2.3  Overview  Idle stats  Frequency stats  Device stats  Tunables

The battery reports a discharge rate of 16.7 W
The estimated remaining time is 1 hours, 25 minutes

Summary: 386.1 wakeups/second, 60.2 GPU ops/seconds, 0.0 VFS ops/sec and 42.9% CPU use

Power est.      Usage      Events/s    Category    Description
3.79 W          2642 rpm      Device      Laptop fan
3.39 W          53.3%        Device      Display backlight
2.63 W          172.9 ms/s    0.00        Timer       process_timeout
2.24 W          142.2 ms/s    17.8        Interrupt   [9] acpi
665 mW         43.6 ms/s     27.5        Process     /usr/lib64/firefox/firefox
237 mW         10.7 ms/s     56.4        Process     /usr/lib64/seamonkey/seamonkey
144 mW         5.7 ms/s      77.2        Interrupt   PS/2 Touchpad / Keyboard / Mouse
119 mW         7.8 ms/s      11.9        Process     /usr/bin/Xorg :0 -background none -verbose -auth /var/run/gdm
91.3 mW        3.7 pkts/s    Device      Network interface: wlan0 (iwlwifi)
84.3 mW        5.5 ms/s     45.9        Timer       tick_sched_timer
77.3 mW        3.3 ms/s     10.1        Process     gkrellm --geometry +1608+70
72.9 mW        4.8 ms/s     20.6        Process     /usr/lib/polkit-1/polkitd --no-debug
58.9 mW        3.9 ms/s     15.0        Process     /usr/lib64/seamonkey/plugin-container /usr/lib64/flash-plugin
51.4 mW        3.4 ms/s     0.00        Interrupt   [1] timer(softirq)
42.3 mW        2.6 ms/s     13.0        Process     xfce4-screenshooter
37.2 mW        2.4 ms/s     58.1        Timer       hrtimer_wakeup
33.0 mW        2.2 ms/s     6.3         Interrupt   [7] sched(softirq)
31.5 mW        60.9 us/s     7.3         kWork      iwl_bg_run_time_calib_work
29.8 mW        2.0 ms/s     41.2        kWork      od_dbt_timer
28.9 mW        1.6 ms/s     1.7         Process     xfce4-panel
25.2 mW        0.9 ms/s     8.6         Process     xfwm4
21.3 mW        1.4 ms/s     0.00        Timer       delayed_work_timer_fn
16.3 mW        1.1 ms/s     0.00        Process     /bin/dbus-daemon --system --address=systemd: --nofork --nopid
13.1 mW        0.9 ms/s     0.5         Process     crond
12.4 mW        0.8 ms/s     0.00        Interrupt   [0] timer/1
12.2 mW        0.8 ms/s     4.3         Interrupt   [6] tasklet(softirq)
12.1 mW        0.8 ms/s     0.05        kWork      disk_events_workfn
12.0 mW        0.8 ms/s     0.00        Interrupt   [0] timer/0
10.0 mW        659.2 us/s    0.4         kWork      kcryptd_crypt
10.0 mW        658.2 us/s    2.1         Process     /usr/sbin/NetworkManager --no-daemon
8.04 mW        528.0 us/s    0.05        Process     powertop
5.76 mW        347.4 us/s    1.6         Process     xchat
5.59 mW        366.9 us/s    0.00        Interrupt   [9] RCU(softirq)
4.75 mW        311.5 us/s    0.00        Process     /usr/sbin/crond -n

```

図2.1 実行中の PowerTOP

PowerTOP を `--html` オプションで実行すると、HTML レポートを生成することもできます。 `htmlfile.html` パラメーターを希望する出力ファイル名に置き換えます。

```
powertop --html=htmlfile.html
```

デフォルトでは、PowerTOP は 20 秒間隔で測定を行います。 `--time` オプションを使うとこれを変更することもできます。

```
powertop --html=htmlfile.html --time=seconds
```

PowerTOP プロジェクトの詳細については、[PowerTOP のホームページ](#)を参照してください。

PowerTOP は `turbostat` ユーティリティーと併用することもできます。 `turbostat` ユーティリティーはレポートングツールで、Intel 64 プロセッサ上のプロセッサポートロジック、周波数、アイドル状態の電力状態、温度、および電力使用量を表示します。 `turbostat` ユーティリティーについての詳細は、`turbostat(8) man` ページまたは [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。

2.3. Diskdevstat と netdevstat

Diskdevstat と netdevstat は、システム上で稼働しているすべてのアプリケーションのディスク活動とネットワーク活動の詳細情報を収集する SystemTap ツールです。これらのツールは、あらゆるアプリケーションで引き起こされる CPU のウェイクアップ回数を秒単位で示す PowerTOP から発想を得ています ([「PowerTOP」](#) を参照)。これらのツールが収集する統計により、小規模な I/O 動作を数多く行なうこ

とで電力を無駄にしているアプリケーションを特定することができますようになります。転送速度のみを測定する他の監視ツールでは、このタイプの使用量を特定できません。

SystemTap のこれらのツールをインストールするには、**root** で以下のコマンドを実行します。

```
yum install tuned-utils-systemtap kernel-debuginfo
```

次のコマンドでツールを実行します。

```
diskdevstat
```

あるいは、以下のコマンドを実行します。

```
netdevstat
```

これらコマンドは両方とも以下のように、最大3つのパラメータを取ります。

```
diskdevstat update_interval total_duration display_histogram
```

```
netdevstat update_interval total_duration display_histogram
```

update_interval

表示が更新される秒単位の間隔。デフォルト: **5**

total_duration

実行完了にかかる秒単位の時間。デフォルト: **86400** (1日)

display_histogram

実行完了時に全収集データによる度数分布図 (柱状グラフ) を作成するかどうかを指定するフラグ。

以下の出力は **PowerTOP** の出力に似ています。以下に KDE 4.2 を稼働している Fedora 10 システム上の長期の **diskdevstat** 実行からのサンプル出力を示します。

PID	UID	DEV	WRITE_CNT	WRITE_MIN	WRITE_MAX	WRITE_AVG	READ_CNT
2789	2903	sda1	854	0.000	120.000	39.836	0
0.000	0.000	0.000	plasma				
15494	0	sda1	0	0.000	0.000	0.000	758
0.000	0.012	0.000	0logwatch				
15520	0	sda1	0	0.000	0.000	0.000	140
0.000	0.009	0.000	perl				
15549	0	sda1	0	0.000	0.000	0.000	140
0.000	0.009	0.000	perl				
15585	0	sda1	0	0.000	0.000	0.000	108
0.001	0.002	0.000	perl				
2573	0	sda1	63	0.033	3600.015	515.226	0
0.000	0.000	0.000	auditd				
15429	0	sda1	0	0.000	0.000	0.000	62
0.009	0.009	0.000	crond				
15379	0	sda1	0	0.000	0.000	0.000	62
0.008	0.008	0.000	crond				
15473	0	sda1	0	0.000	0.000	0.000	62
0.008	0.008	0.000	crond				

15415	0 sda1	0	0.000	0.000	0.000	62
0.008	0.008	0.000 crond				
15433	0 sda1	0	0.000	0.000	0.000	62
0.008	0.008	0.000 crond				
15425	0 sda1	0	0.000	0.000	0.000	62
0.007	0.007	0.000 crond				
15375	0 sda1	0	0.000	0.000	0.000	62
0.008	0.008	0.000 crond				
15477	0 sda1	0	0.000	0.000	0.000	62
0.007	0.007	0.000 crond				
15469	0 sda1	0	0.000	0.000	0.000	62
0.007	0.007	0.000 crond				
15419	0 sda1	0	0.000	0.000	0.000	62
0.008	0.008	0.000 crond				
15481	0 sda1	0	0.000	0.000	0.000	61
0.000	0.001	0.000 crond				
15355	0 sda1	0	0.000	0.000	0.000	37
0.000	0.014	0.001 laptop_mode				
2153	0 sda1	26	0.003	3600.029	1290.730	0
0.000	0.000	0.000 rsyslogd				
15575	0 sda1	0	0.000	0.000	0.000	16
0.000	0.000	0.000 cat				
15581	0 sda1	0	0.000	0.000	0.000	12
0.001	0.002	0.000 perl				
15582	0 sda1	0	0.000	0.000	0.000	12
0.001	0.002	0.000 perl				
15579	0 sda1	0	0.000	0.000	0.000	12
0.000	0.001	0.000 perl				
15580	0 sda1	0	0.000	0.000	0.000	12
0.001	0.001	0.000 perl				
15354	0 sda1	0	0.000	0.000	0.000	12
0.000	0.170	0.014 sh				
15584	0 sda1	0	0.000	0.000	0.000	12
0.001	0.002	0.000 perl				
15548	0 sda1	0	0.000	0.000	0.000	12
0.001	0.014	0.001 perl				
15577	0 sda1	0	0.000	0.000	0.000	12
0.001	0.003	0.000 perl				
15519	0 sda1	0	0.000	0.000	0.000	12
0.001	0.005	0.000 perl				
15578	0 sda1	0	0.000	0.000	0.000	12
0.001	0.001	0.000 perl				
15583	0 sda1	0	0.000	0.000	0.000	12
0.001	0.001	0.000 perl				
15547	0 sda1	0	0.000	0.000	0.000	11
0.000	0.002	0.000 perl				
15576	0 sda1	0	0.000	0.000	0.000	11
0.001	0.001	0.000 perl				
15518	0 sda1	0	0.000	0.000	0.000	11
0.000	0.001	0.000 perl				
15354	0 sda1	0	0.000	0.000	0.000	10
0.053	0.053	0.005 lm_lid.sh				

各コラムの表示内容は、以下のとおりです。

PID

アプリケーションのプロセス ID

UID

アプリケーションの実行元となるユーザー ID

DEV

I/O が発生したデバイス

WRITE_CNT

書き込み動作回数の合計

WRITE_MIN

2 回の連続書き込みに要した最短時間 (秒)

WRITE_MAX

2 回の連続書き込みに要した最長時間 (秒)

WRITE_AVG

2 回の連続書き込みに要した平均時間 (秒)

READ_CNT

読み込み動作回数の合計

READ_MIN

2 回の連続読み込みに要した最短時間 (秒)

READ_MAX

2 回の連続読み込みに要した最長時間 (秒)

READ_AVG

2 回の連続読み込みに要した平均時間 (秒)

COMMAND

プロセスの名前

この例には、非常に目立つアプリケーションが 3 つあります。

PID	UID	DEV	WRITE_CNT	WRITE_MIN	WRITE_MAX	WRITE_AVG	READ_CNT
2789	2903	sda1	854	0.000	120.000	39.836	0
0.000	0.000	0.000	0.000	plasma			
2573	0	sda1	63	0.033	3600.015	515.226	0
0.000	0.000	0.000	0.000	auditd			
2153	0	sda1	26	0.003	3600.029	1290.730	0
0.000	0.000	0.000	0.000	rsyslogd			

これらの 3 つのアプリケーションには、0 以上の **WRITE_CNT** があり、これは測定中になんらかの書き込みを実行したことを意味します。その中でも、**plasma** が他と大差をつけて一番高い数値を示しています。最多の書き込み動作を実行しているため、当然書き込みの平均時間は最短となります。そのため、電力効率の悪いアプリケーションに懸念がある場合には、**Plasma** が調査の最大のターゲットとなります。

strace コマンドと **ltrace** コマンドを使用して、所定のプロセス ID のすべてのシステムコールを追跡することによりアプリケーションをもっと詳しく検査できます。この例では、以下を実行することが可能です。

```
strace -p 2789
```

この例では、**strace** の出力には、ユーザーの KDE アイコンのキャッシュファイルを書き込みのため開き、直後にそのファイルを再び閉じるという動作が 45 秒毎に繰り返されるパターンが含まれていました。これにより、ファイルのメタデータ (特に修正時間) が変更されたため、それに必要な物理的な書き込みがハードディスクに行なわれました。最終的な修正では、アイコンに更新が加えられなかった場合には、こうした不要なコールが発生しないようになりました。

2.4. Battery Life Tool Kit

Red Hat Enterprise Linux 7 では、バッテリーの寿命とパフォーマンスをシミュレートして解析するテストスイートである **BLTK (Battery Life Tool Kit)** が採用されています。BLTK は、このために特定のユーザーグループをシミュレートするタスクセットを実行し、その結果を報告します。特にノートブックのパフォーマンスをテストするために開発された BLTK ですが、**-a** を付けて起動すると、デスクトップコンピュータのパフォーマンスも報告できます。

BLTK を使用すると、実際にマシンを使用しているのと同程度の再現可能な作業負荷を生成できるようになります。例えば、**office** の作業負荷はテキストを書き込み、その中で修正を行います。同じ作業を表計算でも行ないます。BLTK を **PowerTOP** や他の監査用ツール、解析用ツールなどと併用することで、実施した最適化がアイドル状態の時だけでなく、マシンが頻繁に使用されている時にも効果を発揮しているかどうかを検証することができます。全く同じ作業負荷を異なる設定で複数回実行できるため、異なる設定における結果を比較することができます。

以下のコマンドを使用して、BLTK をインストールします。

```
yum install bltk
```

以下のコマンドで、BLTK を実行します。

```
bltk workload options
```

例えば、**idle** の作業負荷を 120 秒間実行するには、以下を実行します。

```
bltk -I -T 120
```

デフォルトで使用できる作業負荷は次の通りです。

-I, --idle

システムはアイドル状態です。他の作業負荷と比較する場合に基準値として使用します。

-R, --reader

ドキュメントを読み込むシミュレートを行います (デフォルトで、**Firefox** を使用)。

-P, --player

CD または DVD ドライブのマルチメディアファイルを見ているシミュレートを行います (デフォルトでは **mplayer** を使用)。

-O, --office

OpenOffice.org スイートを使ったドキュメント編集のシミュレートを行います。

指定できる他のオプションは以下のとおりです。

-a, --ac-ignore

AC 電源が使用可能かどうか無視します (デスクトップで必要)。

-T number_of_seconds, --time number_of_seconds

テストを実行する期間 (秒); **idle** 作業負荷を使ってこのオプションを使用します。

-F filename, --file filename

特定の作業負荷で使用されるファイルを指定します。例えば、CD か DVD ドライブにアクセスする代わりに、**player** の作業負荷で再生するファイルです。

-W application, --prog application

特定の作業負荷で使用されるアプリケーションを指定します。例えば、**reader** の作業負荷で **Firefox** 以外のブラウザを指定します。

BLTK は、数多くの特殊化したオプションをサポートします。詳細情報は **bltk man** ページを参照してください。

BLTK は、生成する結果を **/etc/bltk.conf** 設定ファイルで指定されたディレクトリーに保存します。デフォルトでは **~/bltk/workload.results.number/** です。例えば、**~/bltk/reader.results.002/** ディレクトリーは **reader** の作業負荷の 3 つ目のテスト結果を保持します (ひとつ目のテストは番号なし)。結果はいくつかのテキストファイルに分散されます。これらの結果を読み取りやすい形式に簡略化するには、以下を実行します。

```
bltk_report path_to_results_directory
```

これにより、結果のディレクトリーに **Report** というテキストファイルでその結果が表示されます。代わりにターミナルエミュレーターでその結果を閲覧するには、**-o** オプションを使用します。

```
bltk_report -o path_to_results_directory
```

2.5. Tuned

Tuned は、**udev** を使用して接続されたデバイスを監視し、選択されたプロファイルに従ってシステム設定を動的にチューニングするデーモンです。高スループット、低レイテンシー、省電力などの一般的なユースケース向けの複数の定義済みプロファイルが配布され、各プロファイル向けに定義されたルールを変更し、特定のデバイスのチューニング方法をカスタマイズできます。特定のプロファイルで行われたシステム設定のすべての変更を元に戻すには、別のプロファイルに切り替えるか、**tuned** デーモンを非アクティブ化します。

静的チューニングは、**sysctl** 設定および **sysfs** 設定と **ethtool** などの複数の設定ツールのワンショットアクティベーションから構成されます。また、**Tuned** はシステムコンポーネントの使用を監視し、その監視情報に基づいてシステム設定を動的にチューニングします。動的なチューニングでは、該当するシステムの稼働時間中に各種システムコンポーネントを異なる方法で使うことが考慮されます。たとえば、ハードドライブは起動時とログイン時に頻繁に使用されますが、ユーザーが主に Web ブラウザーや電子メールクライアントなどのアプリケーションを使用する場合はほとんど使用されません。同様に、CPU とネットワークデバイスは、異なるタイミングで異なる方法で使われます。**Tuned** はこれらのコンポーネントの動作を監視し、それらの使用の変化に反応します。


注記

動的チューニングは Red Hat Enterprise Linux ではグローバルで無効であり、`/etc/tuned/tuned-main.conf` ファイルを編集し、`dynamic_tuning` フラグを **1** に変更することによって有効にできます。

実際の例として、標準的なオフィスにあるワークステーションを考えます。通常、イーサネットのネットワークインターフェースはほとんど使用されず、ほんの数件の電子メールがたまに送受信されるか、ウェブページが数ページ読み込まれる程度だとします。このような負荷の場合、ネットワークインターフェースにデフォルト設定のように常に最高速度で動作する必要はありません。**Tuned** には、ネットワークデバイスを監視してチューニングを行なうプラグインがあり、ネットワークインターフェースの利用率が低くなると、それを検出して自動的にそのインターフェースの速度を低下させて電力の使用を削減することができます。たとえば、DVD イメージをダウンロードしたり、大きい添付ファイルが付いた電子メールを開いたりしたためインターフェースのアクティビティーが長い時間にわたって増加すると、**tuned** はこれを検出して、アクティビティレベルが非常に高い間に最良のパフォーマンスを提供できるようにインターフェース速度を最大に設定します。この原則は CPU およびハードディスク向けの他のプラグインにも使用されています。

2.5.1. プラグイン

一般的に、**tuned** では、**監視プラグイン**と**チューニングプラグイン**の2つの種類のプラグインを使用します。監視プラグインは、稼働中のシステムから情報を取得するために使用されます。現時点では、以下の監視プラグインが実装されています。

disk

デバイスおよび測定間隔ごとのディスク負荷 (IO 操作の数) を取得します。

net

ネットワークカードおよび測定間隔ごとのネットワーク負荷 (転送済みパケットの数) を取得します。

load

CPU および測定間隔ごとの CPU 負荷を取得します。

監視プラグインの出力は、動的チューニング向けチューニングプラグインによって使用できます。現在実装されている動的チューニングアルゴリズムは、パフォーマンスと省電力のバランスを取ろうとし、パフォーマンスプロファイルで無効になります (個別プラグインの動的チューニングは **tuned** プロファイルで有効または無効にできます)。監視プラグインは、有効ないずれかのチューニングプラグインでメトリクスが必要な場合に必ず自動的にインスタンス化されます。2つのチューニングプラグインで同じデータが必要な場合は、監視プラグインのインスタンスが1つだけ作成され、データが共有されます。

各チューニングプラグインにより、個別サブシステムがチューニングされ、**tuned** プロファイルから入力される複数のパラメーターが取得されます。各サブシステムには、チューニングプラグインの個別インスタンスで処理される複数のデバイス (複数の CPU やネットワークカードなど) を含めることができます。また、個別デバイスの特定の設定もサポートされます。提供されたプロファイルでは、個別サブシステムのすべてのデバイスに一致するワイルドカードが使用されます (このプロファイルの変更方法の詳細については、[「カスタムプロファイル」](#)を参照)。これにより、プラグインは必要な目標 (選択されたプロファイル) に従ってこれらのサブシステムをチューニングできるようになります。ユーザーが行うのは、正しい **tuned** プロファイルを選択するだけです (プロファイルの選択方法の詳細または提供されたプロファイルのリストについては、[「Tuned」](#)と[「Tuned-adm」](#)を参照)。現時点では、以下のチューニングプラグインが実装されています (動的チューニングはこれらの一部のプラグインのみで実装され、プラグインでサポートされたパラメーターもリストされます)。

cpu

CPU ガバナーを、**governor** パラメーターで指定された値に設定し、CPU 負荷に応じて PM QoS CPU DMA レイテンシーを動的に変更します。CPU 負荷が **load_threshold** パラメーターで指定された値よりも小さい場合、レイテンシーは **latency_high** パラメーターで指定された値に設定されます。それ以外の場合は、**latency_low** で指定された値に設定されます。また、レイテンシーは、特定の値に強制的に設定できます (動的に変更しない場合)。これは、**force_latency** パラメーターを必要なレイテンシー値に設定することにより実現できます。

eeepc_she

CPU 負荷に応じて FSB 速度を動的に設定します。この機能は、一部のノートブックに存在し、Asus Super Hybrid Engine と呼ばれます。CPU 負荷が **load_threshold_powersave** パラメーターで指定された値以下である場合は、プラグインによって FSB 速度が **she_powersave** パラメーターで指定された値に設定されます (FSB 周波数および対応する値の詳細については、カーネルのドキュメントを参照)。CPU 負荷が **load_threshold_normal** パラメーターで指定された値以上である場合は、FSB 速度が **she_normal** パラメーターで指定された値に設定されます。静的チューニングはサポートされず、プラグインは透過的に無効になります (この機能のハードウェアサポートが検出されない場合)。

net

wake-on-lan を **wake_on_lan** パラメーターで指定された値に設定します (**ethtool** ユーティリティと同じ構文を使用します)。また、インターフェースの使用状況に応じてインターフェース速度が自動的に変更されます。

sysctl

プラグインパラメーターで指定されたさまざまな **sysctl** 設定を指定します。構文は **name=value** となります。ここで、**name** は **sysctl** ツールにより提供されたものと同じ名前になります。このプラグインは、他のプラグインで指定できない設定を変更する必要がある場合に使用します (ただし、設定を特定のプラグインで指定できる場合は、そのプラグインを使用することが推奨されます)。

usb

USB デバイスの autosuspend タイムアウトを **autosuspend** パラメーターで指定された値に設定します。値が 0 の場合は、autosuspend が無効になります。

vm

transparent_hugepages パラメーターのブール値に応じて透過的な大規模ページを有効または無効にします。

audio

音声コーデックの autosuspend タイムアウトを **timeout** パラメーターで指定された値に設定します。現時点では、**snd_hda_intel** と **snd_ac97_codec** がサポートされています。値が 0 の場合は、autosuspend が無効になります。また、ブール値パラメーター **reset_controller** を **true** に設定することにより、コントローラーを強制的にリセットすることもできます。

disk

エレベーターを **elevator** パラメーターで指定された値に設定します。また、ALPM を **alpm** パラメーターで指定された値 ([「Aggressive Link Power Management」](#)を参照)、ASPM を **aspm** パラメーターで指定された値 ([「Active-State Power Management」](#)を参照)、スケジューラークォンタムを **scheduler_quantum** パラメーターで指定された値、ディスク spindown タイ

ムアウトを **spindown** パラメーターで指定された値、ディスク readahead を **readahead** パラメーターで指定された値に設定し、現在のディスク readahead 値を **readahead_multiply** パラメーターで指定された定数で乗算できます。さらに、このプラグインにより、現在のドライブ使用状況に応じてドライブの高度な電力管理設定と spindown タイムアウト設定が動的に変更されます。動的チューニングは、ブール値パラメーター **dynamic** により制御でき、デフォルトで有効になります。

mounts

disable_barriers パラメーターのブール値に応じてマウントのバリアを有効または無効にします。

script

このプラグインは、プロファイルがロードまたはアンロードされたときに実行する外部スクリプトの実行に使用されます。スクリプトは、**start** または **stop** のいずれかの引数 (スクリプトがプロファイルのロード時またはアンロード時に呼び出されるかによって決まります) によって呼び出されます。スクリプトファイル名は、**script** パラメーターによって指定できます。スクリプトで停止アクションを正しく実装し、変更したすべての設定を起動アクション中に元に戻す必要があることに注意してください。このような操作を行わないと、ロールバックが機能しません。ユーザーの利便性のために、**functions** Bash ヘルパースクリプトがデフォルトでインストールされ、スクリプトで定義されたさまざまな機能をインポートして使用できます。この機能は、主に後方互換性を維持するために提供されます。この機能は最終手段として使用し、必要な設定を指定できる他のプラグインが存在する場合は、そのプラグインを使用することが推奨されます。

sysfs

プラグインパラメーターで指定されたさまざまな **sysfs** 設定を指定します。構文は **name=value** となります。ここで、**name** は、使用する **sysfs** パスです。このプラグインは、他のプラグインで指定できない設定を変更する必要がある場合に使用します (設定を特定のプラグインで指定できる場合は、そのプラグインを使用することが推奨されます)。

video

ビデオカードのさまざまな省電力レベルを設定します (現時点では、Radeon カードのみがサポートされます)。省電力レベルは **radeon_powersave** パラメーターを使用して指定できます。サポートされる値は **default**、**auto**、**low**、**mid**、**high**、および **dynpm** です。詳細については、http://www.x.org/wiki/RadeonFeature#KMS_Power_Management_Options を参照してください。このプラグインは試験目的で提供され、今後のリリースで変更されることがあることに注意してください。

2.5.2. インストールと使用方法

tuned パッケージをインストールするには、**root** で以下のコマンドを実行します。

```
yum install tuned
```

tuned パッケージのインストールでは、お使いのシステムに最適なプロファイルも事前に設定されます。現時点では、以下のカスタマイズ可能なルールに従ってデフォルトのプロファイルが選択されます。

throughput-performance

これは、コンピュータノードとして動作する Fedora オペレーティングシステムで事前に選択されます。このようなシステムの目的は、スループットパフォーマンスの最大化です。

virtual-guest

これは、仮想マシンで事前に選択されます。この目的はパフォーマンスの最大化です。パフォーマンスの最大化に興味がない場合は、**balanced** または **powersave** プロファイルに変更できます (下記参照)。

balanced

これは、他のすべてのケースで事前に選択されます。この目的は、パフォーマンスと電力消費の調和です。

tuned を起動するには、**root** で以下のコマンドを実行します。

```
systemctl start tuned
```

マシンを起動する度に **tuned** を有効にするには、以下のコマンドを実行します。

```
systemctl enable tuned
```

プロファイルの選択などの他の **tuned** の制御の場合は、以下のコマンドを実行します。

```
tuned -adm
```

このコマンドを使用するには、**tuned** サービスが実行されている必要があります。

インストールされた利用可能なプロファイルを参照するには、以下のコマンドを実行します。

```
tuned -adm list
```

現在アクティブなプロファイルを参照するには、以下のコマンドを実行します。

```
tuned -adm active
```

プロファイルを選択またはアクティブ化するには、以下のコマンドを実行します。

```
tuned -adm profile profile
```

例えば：

```
tuned -adm profile powersave
```

試験目的で提供された機能として、複数のプロファイルを一度に選択することができます。**tuned** アプリケーションは、ロード中にそれらのプロファイルをマージしようとします。競合が発生した場合は、最後に指定されたプロファイルの設定が優先されます。これは自動的に行われ、使用されるパラメーターの組み合わせが適切であるかどうかはチェックされません。あまり深く考えずにこの機能を使用すると、一部のパラメーターが反対の方向でチューニングされ、生産性が上がらないことがあります。このような状況の例として、**throughput-performance** プロファイルでディスクに対して **high** スループットを設定し、同時に **spindown-disk** プロファイルでディスクスピンドアウンを **low** 値に設定することが挙げられます。以下の例では、仮想マシンでの実行でパフォーマンスを最大化するようシステムが最適化され、同時に、低消費電力を実現するようシステムがチューニングされます (低消費電力が最優先である場合)。

```
tuned -adm profile virtual-guest powersave
```

既存のプロファイルを変更したり、インストール中に使用されたのと同じロジックを使用したりせずに **tuned** がシステムに最適なプロファイルを提示するようにするには、以下のコマンドを実行します。

tuned -adm recommend

Tuned には、手動で実行する場合に使用できる追加オプションがあります。ただし、このオプションは推奨されず、主にデバッグ目的で提供されています。利用可能なオプションは、以下のコマンドを使用して表示できます。

tuned --help**2.5.3. カスタムプロファイル**

ディストリビューション固有のプロファイルは、`/usr/lib/tuned` ディレクトリーに格納されます。各プロファイルには独自のディレクトリーがあります。プロファイルは **tuned.conf** という名前の主要設定ファイルとオプションのヘルプスクリプトなどの他のファイルから構成されます。`/usr/lib/tuned` 内のプロファイルは変更しないでください。プロファイルのカスタマイズする必要がある場合は、プロファイルディレクトリーを `/etc/tuned` ディレクトリーにコピーします。同じ名前のプロファイルが2つある場合は、`/etc/tuned` のプロファイルが優先されます。また、興味があるプロファイルを含む `/etc/tuned` ディレクトリーで独自のプロファイルを作成し、必要なパラメーターのみを変更または上書きすることもできます。

tuned.conf ファイルには、複数のセクションが含まれます。**[main]** セクションは1つだけ存在します。他のセクションはプラグインインスタンス向けの設定です。**[main]** セクションを含むすべてのセクションはオプションです。コメントもサポートされます。ハッシュ (#) で始まる行はコメントです。

[main] セクションには、以下のオプションがあります。

include=profile

指定されたプロファイルがインクルードされます。たとえば、**include=powersave** の場合は、**powersave** プロファイルがインクルードされます。

プラグインインスタンスが記述されるセクションは、以下のように書式化されます。

```
[NAME]
type=TYPE
devices=DEVICES
```

NAME は、ログで使用されるプラグインインスタンスの名前であり、任意の文字列です。**TYPE** は、チューニングプラグインのタイプです。チューニングプラグインのリストと説明については、「[プラグイン](#)」を参照してください。**DEVICES** は、このプラグインインスタンスが処理するデバイスのリストです。**devices** 行には、リスト、ワイルドカード (*), および否定 (!) を含めることができます。また、ルールを組み合わせることもできます。**devices** 行がない場合は、**TYPE** のシステムに存在する、または後で接続されたすべてのデバイスがプラグインインスタンスによって処理されます。これは、**devices=*** を使用したときと同様です。プラグインのインスタンスが指定されない場合、プラグインは無効になりません。プラグインがさらに多くのオプションをサポートする場合は、プラグインセクションでそれらのプラグインを指定することもできます。オプションが指定されないと、デフォルト値が使用されます (インクルードされたプラグインで以前に指定されていない場合)。プラグインオプションのリストについては、「[プラグイン](#)」を参照してください。

例2.1 プラグインインスタンスの定義

以下の例では、**sda** や **sdb** などの **sd** で始まるすべての候補に一致し、それらの候補に対するバリアは無効になりません。

```
[data_disk]
type=disk
```

```
devices=sd*
disable_barriers=false
```

以下の例では、**sda1** と **sda2** を除くすべての候補に一致します。

```
[data_disk]
type=disk
devices=!sda1, !sda2
disable_barriers=false
```

プラグインインスタンスのカスタム名を必要とせず、設定ファイルでインスタンスの定義が1つしかない場合、Tuned は以下の短い構文をサポートします。

```
[TYPE]
devices=DEVICES
```

この場合は、**type** 行を省略することができます。タイプと同様に、インスタンスは名前で参照されます。上記の例は、以下のように書き換えることができます。

```
[disk]
devices=sdb*
disable_barriers=false
```

include オプションを使用して同じセクションが複数回指定された場合は、設定がマージされます。競合のため、設定をマージできない場合は、競合がある以前の設定よりも競合がある最後の定義が優先されます。場合によっては、以前に定義された内容がわからないことがあります。このような場合は、**replace** ブール値オプションを使用して **true** に設定できます。これにより、同じ名前の以前の定義がすべて上書きされ、マージは行われません。

また、**enabled=false** オプションを指定してプラグインを無効にすることもできます。これは、インスタンスが定義されない場合と同じ効果を持ちます。**include** オプションから以前の定義を再定義し、カスタムプロファイルでプラグインをアクティブにしない場合は、プラグインを無効にすると便利です。

ほとんどの場合、デバイスは1つのプラグインインスタンスで処理できます。デバイスが複数のインスタンス定義に一致する場合は、エラーが報告されます。

以下に、**balanced** プロファイルに基づき、すべてのデバイスの ALPM が最大の省電力に設定されるように拡張されたカスタムプロファイルの例を示します。

```
[main]
include=balanced

[disk]
alpm=min_power
```

2.5.4. Tuned-adm

システムを詳細に監査、分析することは、非常に時間のかかる作業であり、数ワットの節電のためだけに行う価値はないかもしれません。今までは、こうした煩雑な作業の代わりになる方法は、単にデフォルトを使用することだけでした。そのため、Red Hat Enterprise Linux 7 では、こうした両極端な方法の代替案となるプロファイルを複数用意し、特定の使用例に対応しています。さらには、コマンドラインで複数のプロファイルを簡単に切り替えることができる **tuned-adm** ツールも提供しています。Red Hat Enterprise Linux 7 には、一般的な使用例に則した定義済みのプロファイルが同梱されているので、**tuned-adm** コマンドで選択してアクティベートするだけで使用できるようになります。ただし、プロファイルをユーザー自

身で作成したり、修正を行なうほか、削除することも可能です。

利用可能な全プロファイルを一覧表示して、現在アクティブなプロファイルを特定するには、以下を実行します。

```
tuned-adm list
```

現在アクティブなプロファイルだけを表示する場合は、以下を実行します。

```
tuned-adm active
```

別のプロファイルに切り替える場合は、以下を実行します。

```
tuned-adm profile profile_name
```

例を示します。

```
tuned-adm profile server-powersave
```

すべてのチューニングを無効にする場合は、以下を実行します。

```
tuned-adm off
```

以下に、基本パッケージでインストールされるプロファイルをリストします。

balanced

デフォルトの省電力プロファイル。パフォーマンスと電力消費のバランスを取ることが目的です。可能な限り、自動スケーリングと自動チューニングを使用しようとしています。ほとんどの負荷で良い結果をもたらします。唯一の欠点はレイテンシーが増加することです。現在の **tuned** リリースでは、CPU、ディスク、音声、および動画のプラグインが有効になり、**ondemand** ガバナーがアクティブ化されます。**radeon_powersave** は **auto** に設定されます。

powersave

省電力パフォーマンスを最大化するプロファイル。実際の電力消費を最小化するためにパフォーマンスを調整できます。現在の **tuned** リリースでは、USB 自動サスペンド、WiFi 省電力、および SATA ホストアダプター向けの ALPM 省電力が有効になります ([「Aggressive Link Power Management」](#)を参照)。また、ウェイクアップ率が低いシステムのマルチコア省電力がスケジュールされ、**ondemand** ガバナーがアクティブ化されます。さらに、AC97 音声省電力と、システムに応じて HDA-Intel 省電力 (10 秒のタイムアウト) が有効になります。KMS が有効なサポート対象 Radeon グラフィックカードがシステムに搭載されている場合は、自動省電力に設定されます。Asus Eee PC では、動的な Super Hybrid Engine が有効になります。



注記

powersave プロファイルは、必ずしも最も効率的ではありません。定義された量の作業を行う場合 (たとえば、動画ファイルをトランスコードする必要がある場合) を考えてください。トランスコードがフルパワーで実行される場合に、マシンが少ない電力を消費することがあります。これは、タスクがすぐに完了し、マシンがアイドル状態になり、非常に効率的な省電力モードに自動的に切り替わることがあるためです。その一方で、調整されたマシンでファイルをトランスコードすると、マシンはトランスコード中に少ない電力を消費しますが、処理に時間がかかり、全体的な消費電力は高くなる場合があります。このため、一般的に **balanced** プロファイルが優れたオプションになる場合があります。

throughput-performance

高スループットに最適化されたサーバープロファイル。省電力メカニズムが無効になり、ディスクとネットワーク IO のスループットパフォーマンスを向上させる sysctl 設定が有効になり、**deadline** スケジューラーに切り替わります。CPU ガバナーは **performance** に設定されます。

latency-performance

低レイテンシーに最適化されたサーバープロファイル。省電力メカニズムが無効になり、レイテンシーを向上させる sysctl 設定が有効になります。CPU ガバナーは **performance** に設定され、CPU は低い C 状態にロックされます (PM QoS を使用)。

network-latency

低レイテンシーネットワークチューニング向けプロファイル。**latency-performance** プロファイルに基づきます。また、透過的な巨大ページと NUMA 調整が無効になり、複数の他のネットワーク関連の sysctl パラメーターがチューニングされます。

network-throughput

スループットネットワークチューニング向けプロファイル。**throughput-performance** プロファイルに基づきます。また、カーネルネットワークバッファが増加されます。

virtual-guest

enterprise-storage プロファイルに基づく仮想ゲスト向けプロファイル。仮想メモリーの swappiness の減少やディスクの readahead 値の増加などが行われます。ディスクバリアは無効になりません。

virtual-host

enterprise-storage プロファイルに基づく仮想ホスト向けプロファイル。仮想メモリーの swappiness の減少、ディスクの readahead 値の増加、ダーティーページのアグレッシブライトバックの有効化などが行われます。

sap

SAP ソフトウェアのパフォーマンスを最大化するように最適化されたプロファイルです。**enterprise-storage** プロファイルをベースにしています。sap プロファイルは、共有メモリー、セマフォ、プロセスが所有するメモリーマップの最大数に関する sysctl 設定も調整します。

desktop

balanced プロファイルに基づく、デスクトップに最適化されたプロファイル。対話型アプリケーションの応答を向上させるためにスケジューラーオートグループが有効になります。

追加の定義済みプロファイルは、**Optional** チャンネルで利用可能な *tuned-profiles-compat* パッケージでインストールできます。これらのプロファイルは、後方互換性を維持するために提供され、開発は行われていません。基本パッケージの一般的なプロファイルを使用すると、ほとんどの場合、同等のことやそれ以外のことも行えます。特別な理由がない限り、基本パッケージのプロファイルを使用することが推奨されます。互換プロファイルは以下のとおりです。

default

これは利用可能なプロファイルの中で省電力への影響度が最も低く、**tuned** の CPU プラグインとディスクプラグインのみが有効になります。

desktop-powersave

デスクトップシステム向けの節電プロファイルです。SATA ホストアダプター用の ALPM 節電 ([「Aggressive Link Power Management」](#) 参照) と **tuned** の CPU、イーサネット、およびディスクプラグインを有効にします。

server-powersave

サーバーシステム向けの省電力プロファイル。SATA ホストアダプター用の ALPM 省電力が有効になり、**tuned** の CPU プラグインとディスクプラグインがアクティブ化されます。

laptop-ac-powersave

AC 電源で稼働するノートパソコン向け省電力プロファイル (影響度は中)。SATA ホストアダプター用の ALPM 省電力、WiFi 省電力、および **tuned** の CPU、イーサネット、ディスクのプラグインが有効になります。

laptop-battery-powersave

バッテリーで稼働するラップトップ向け省電力プロファイル (影響度は大)。現在の **tuned** 実装では、**powersave** プロファイルのエイリアスになります。

spindown-disk

スピンドウン時間を最大化する、従来の HDD が搭載されたマシン向け省電力プロファイル。**tuned** 省電力メカニズム、USB 自動サスペンド、および Bluetooth が無効になり、Wi-Fi 省電力が有効になり、ログ同期が無効になり、ディスクライトバック時間が増加し、ディスクの swappiness が減少します。すべてのパーティションは、**noatime** オプションを使用して再マウントされます。

enterprise-storage

I/O スループットを最大化する、エンタープライズ級のストレージ向けサーバープロファイル。**throughput-performance** プロファイルと同じ設定がアクティブ化され、**readahead** 設定値が乗算され、ルートパーティションと起動パーティション以外のパーティションでバリアが無効になります。

2.5.5. Powertop2tuned

powertop2tuned ユーティリティは、**PowerTOP** の提案からカスタム **tuned** プロファイルを作成できるツールです。**PowerTOP** の詳細については、[「PowerTOP」](#) を参照してください。

powertop2tuned アプリケーションをインストールするには、以下のコマンドを root で実行します。

```
yum install tuned-utils
```

カスタムプロファイルを作成するには、以下のコマンド root で実行します。

```
powertop2tuned new_profile_name
```

デフォルトでは、現在選択されている **tuned** プロファイルに基いて **/etc/tuned** ディレクトリー内にプロファイルが作成されます。安全上の理由から、すべての **PowerTOP** チューニングは最初に新しいプロファイルで無効になっています。これらのチューニングを有効にするには、**/etc/tuned/profile/tuned.conf** で興味があるチューニングをコメント解除します。--**enable** または **-e** オプションを使用して、有効な **PowerTOP** により提示されたほとんどのチューニングで新しいプロファイルを生成できます。USB 自動サスペンドなどの一部の危険なチューニングは引き続き無効になります。これらのチューニングが必要な場合は、手動でコメント解除する必要があります。デフォルトでは、新しいプロファイルはアクティブ化されません。アクティブ化するには、以下のコマンドを実行します。

```
tuned-adm profile new_profile_name
```

powertop2tuned がサポートするオプションの完全な一覧を表示するには、以下のコマンドを実行します。

```
powertop2tuned --help
```

2.6. UPower

Red Hat Enterprise Linux 6 では、**DeviceKit-power** は、**HAL** の一部である電力管理機能と Red Hat Enterprise Linux の以前のリリースの **GNOME Power Manager** の一部である電力管理機能を引き継ぎました ([「GNOME の電源管理」](#) も参照)。Red Hat Enterprise Linux 7 では、**DeviceKit-power** は **UPower** という名前に変更されました。**UPower** は、デーモン、API、および一連のコマンドラインツールを提供します。物理デバイスかどうかに関係なく、システム上の各電源はデバイスとして表されます。たとえば、ノートパソコンのバッテリーと AC 電源は両方ともデバイスとして表されます。

コマンドラインツールにアクセスするには、**upower** コマンドと以下のオプションを使用します。

--enumerate, -e

システム上の電源デバイス用のオブジェクトパスを表示します。例えば以下のとおりです。

```
/org/freedesktop/UPower/devices/line_power_AC
/org/freedesktop/UPower/devices/battery_BAT0
```

--dump, -d

システム上の全ての電源デバイス用のパラメータを表示します。

--wakeups, -w

システムの CPU のウェイクアップを表示します。

--monitor, -m

AC 電源の接続や切断、あるいはバッテリーの低下などの電源デバイスの変化についてシステムを監視します。システムの監視を止めるには、**Ctrl+C** を押します。

--monitor-detail

AC 電源の接続や切断、あるいはバッテリーの低下などの電源デバイスの変化についてシステムを監視します。**--monitor-detail** オプションでは、**--monitor** オプションよりも詳細を提供します。システムの監視を止めるには、**Ctrl+C** を押します。

```
--show-info object_path, -i object_path
```

特定のオブジェクトパスに利用可能なすべての情報が表示されます。たとえば、オブジェクトパス `/org/freedesktop/UPower/devices/battery_BAT0` で表されるシステムのバッテリーに関する情報を取得するには、以下のコマンドを実行します。

```
upower -i /org/freedesktop/UPower/devices/battery_BAT0
```

2.7. GNOME の電源管理

GNOME Power Manager は、GNOME デスクトップの一部としてインストールされるデーモンです。Red Hat Enterprise Linux の以前のバージョンで **GNOME Power Manager** が提供した電力管理機能の大部分は、Red Hat Enterprise Linux 6 で **DeviceKit-power** の一部となり、Red Hat Enterprise Linux 7 では **UPower** という名前に変更されました（「[UPower](#)」を参照）。ただし、**GNOME Power Manager** はその機能のフロントエンドとして残ります。システムトレイのアプレットを介して **GNOME Power Manager** は、バッテリーから AC 電源への切り替えなど、システムの電源状態の変化を通知します。また、バッテリーの状態を報告し、バッテリーの電力が低くなると警告を出します。

2.8. 他の監査ツール

Red Hat Enterprise Linux7 は、システムの監査と分析を実行する複数のツールを提供します。それらのほとんどは、すでに発見したものを検証する場合や特定の部分の詳細情報が必要な場合に補助の情報源として使用できます。これらのツールの多くはパフォーマンスチューニングにも使用されます。以下に、これらのツールを示します。

vmstat

vmstat はプロセス、メモリ、ページング、ブロック I/O、トラップ、および CPU 活動について詳細情報を提供します。システム全体で実行している動作やビジーな部分を詳しく見るために使用します。

iostat

iostat は **vmstat** と似ていますが、ブロックデバイスの I/O 専用です。詳細な出力と統計も提供します。

blktrace

blktrace は、非常に詳細に渡るブロック I/O のトレースプログラムです。情報をアプリケーションに関連した 1 つずつのブロックに分割します。**diskdevstat** と併せて使用すると大変役立ちます。

第3章 中核となるインフラストラクチャとメカニズム



重要

本章で解説している `cpupower` コマンドを使用する場合は、`cpupowerutils` パッケージがインストールされていることを確認してください。

3.1. CPU のアイドル状態

x86 アーキテクチャの CPU は、CPU の一部を停止させる設定や低いパフォーマンスで実行させる設定など、様々な状態に対応します。こうした状態は **C 状態** と呼ばれ、使用されていない CPU を部分的に停止させることで節電を可能にしています。C 状態は番号付けされ、C0 から始まり数値が増えていきます。大きな数字ほど、CPU の機能性は低く節電率は高くなります。特定の番号が付いた C 状態は、プロセッサ間でさほど違いませんが、特定の機能に関しては正確にはプロセッサファミリー間で異なる場合があります。C 状態 0 から 3 は以下のように定義されています。

C0

稼働中、または実行中の状態。この状態では、CPU は完全に動作中であり、アイドル状態の部分はありません。

C1, 停止

プロセッサが何の指示も実行していない状態ですが、一般的には電力が低い状態でもありません。CPU は実質的に遅延なく処理を継続できます。C 状態を提供するプロセッサはすべて、この状態に対応できなければなりません。Pentium 4 のプロセッサは、実際には電力消費が低い状態の C1E と呼ばれる拡張型 C1 状態に対応します。

C2, クロック停止

このプロセッサのクロックが停止している状態ですが、そのレジスタとキャッシュの状態は完全な状態で保持しているため、クロックを再開させると直ちに処理を再開することができます。オプションの状態になります。

C3, スリープ

プロセッサが実際にスリープ状態に入り、キャッシュの更新をする必要がない状態です。この状態から復帰するには、C2 状態からの復帰に比べ、かなり長い時間がかかります。この状態もオプションになります。

利用可能なアイドル状態および `CPUIidle` ドライバーの統計値を表示させるには、次のコマンドを実行します。

```
cpupower idle-info
```

Nehalem マイクロアーキテクチャを搭載する近年の Intel CPU の特徴は、新しい C 状態である C6 です。これにより、CPU の電圧供給をゼロにまで下げることが可能ですが、通常は 80% から 90% 電力消費量を低減します。Red Hat Enterprise Linux 7 のカーネルには、この新しい C 状態に対する最適化が含まれています。

3.2. CPUfreq ガバナーの使用

ご使用のシステムで電力消費と発生熱量を低減する効果的な方法の1つは、CPUfreq を使用することです。CPU 速度スケールとも呼ばれる CPUfreq により、プロセッサのクロック速度をオンザフライで調節できます。これにより、システムは減速したクロック速度で稼働し、節電します。CPUfreq ガバナーが、クロック速度の変更、周波数を変換する時期の決定といった周波数の変換に関する規則を定義します。

ガバナーは、システムの CPU の電力特性を定義し、これは結果的に CPU のパフォーマンスに影響を与えます。ガバナーには作業負荷に関してそれぞれ特有の動作、目的、および適合性があります。このセクションでは、CPUfreq ガバナーの選択および設定方法、各ガバナーの特性、およびガバナーに適している作業負荷の種類について説明します。

3.2.1. CPUfreq ガバナーのタイプ

このセクションでは、Red Hat Enterprise Linux 7 で利用できる様々な種類の CPUfreq ガバナーについて説明しています。

cpufreq_performance

Performance ガバナーは、CPU が最高クロック周波数を使用するように強制します。この周波数は静的に設定され、変化しないため、このガバナーでは、**節電する利点はありません**。このガバナーは、何時間にも渡るような作業負荷が大きい時だけ、しかも CPU がアイドル状態になることがほとんどない（もしくはまったくならない）時のみに適しています。

cpufreq_powersave

一方、Powersave ガバナーは、CPU が最低クロック周波数を使用するように強制します。この周波数は静的に設定され変化しないため、このガバナーでは最大の節電を実現しますが、**CPU パフォーマンスが一番低く**なってしまいます。

しかし「節電 (Powersave)」という用語は時に誤解を招きます。全負荷で遅い CPU は (原則として)、負荷がない高速の CPU よりも多くの電力を消費します。そのため、低活動が予期できる時には Powersave ガバナーを使用するよう CPU を設定することが推奨されますが、この期間中に予期しない高負荷が発生するとシステムは実際にはより多くの電力を消費することがあります。

Powersave ガバナーは簡単にいうと、CPU にとっては「節電」よりも「スピードリミッター」の意味を持ちます。これは、過熱が問題となる恐れがあるシステムや環境で最も役立ちます。

cpufreq_ondemand

Ondemand ガバナーは動的なガバナーです。システム負荷が大きい時は、CPU は最高クロック周波数を実現し、システムがアイドル状態の時には、CPU は最低クロック周波数を実現します。これにより、システム負荷に対してシステムは電力消費量を適宜調節できますが、そうすることで **周波数変換の間の遅延**が発生してしまいます。そのため、システムがアイドル状態と高負荷の間で頻繁に替わりすぎると、遅延により、Ondemand ガバナーが実現できるパフォーマンスおよびまたは節電の利点が少なくなる恐れがあります。

ほとんどのシステムでは、Ondemand ガバナーは熱の放出、消費電力、パフォーマンス、および管理のしやすさの間で、最良の妥協策を提供します。1 日の中で特定の時間帯にのみシステムがビジーになる場合は、Ondemand ガバナーはそれ以上介入せずに、負荷に応じて最高周波数と最低周波数の間で自動的に切り替わります。

cpufreq_userspace

Userspace ガバナーを使用すると、ユーザースペースプログラム (または、root で実行しているいずれのプロセス) が周波数を設定できます。Userspace ガバナーは、すべてのガバナーの中で最もカスタマイズ可能であり、設定によってはご使用のシステムでパフォーマンスと電力消費のバランスを最適化できます。

cpufreq_conservative

Ondemand ガバナーと同様に、Conservative ガバナーも使用量に応じてクロック周波数を調節します (Ondemand ガバナーと同様です)。ただし、Ondemand ガバナーがより積極的にクロック周波数を調節するのに対し (最高周波数から最低周波数、そして最高周波数に戻る)、Conservative ガバナーはもっとゆっくりと調節を行います。

これが意味しているのは、Conservative ガバナーは単に最高と最低の周波数を選択するのではなく、負荷に対して適切と判断するクロック周波数に合わせるということです。これは電力消費に著しく貢献する可能性があります。Ondemand ガバナーよりも長い遅延で行います。



注記

cron ジョブを使用してガバナーを有効にできます。これにより、1 日のある時間帯にあるガバナーを自動的に設定することができます。そのため、アイドル状態 (例えば終業後) の時には、低周波数のガバナーを指定し、高負荷となる時間帯には高周波数に戻るよう設定できます。

特定のガバナーを有効にする方法については、[「CPUfreq の設定」](#)を参照してください。

3.2.2. CPUfreq の設定

すべての CPUfreq ドライバーは、kernel-tools パッケージの一部としてビルドされ、自動的に選択されます。したがって、ガバナーを選択するだけで CPUfreq をセットアップできます。

以下のコマンドを実行すると、特定の CPU に使用できるガバナーを表示できます。

```
cpupower frequency-info --governors
```

以下のコマンドを実行すると、すべての CPU に対してこれらのいずれかのガバナーを有効にできます。

```
cpupower frequency-set --governor [governor]
```

特定のコアに対してのみガバナーを有効にするには、CPU メンバーの範囲またはカンマ区切りリストとともに **-c** を使用します。たとえば、CPU 1~3 および 5 の Userspace ガバナーを有効にするには、以下のコマンドを実行します。

```
cpupower -c 1-3,5 frequency-set --governor cpufreq_userspace
```

3.2.3. CPUfreq ポリシーおよび速度のチューニング

該当する CPUfreq ガバナーを選択すると、**cpupower frequency-info** コマンドで CPU 速度とポリシー情報を表示させることができますようになります。さらに、**cpupower frequency-set** のオプションを使うと各 CPU の速度を調整することができます。

cpupower frequency-info には、以下のようなオプションを使用することができます。

- ✧ **--freq** — CPUfreq コアに準じて現在の CPU の速度を KHz 単位で表示します。
- ✧ **--hwfreq** — ハードウェアに準じて現在の CPU の速度を KHz 単位で表示します (root による実行のみ可)。
- ✧ **--driver** — この CPU で周波数の設定に使用している CPUfreq ドライバーを表示します。

- ※ **--governors** — このカーネルで使用できる CPUfreq ガバナーを表示します。このファイルには表示されていない CPUfreq ガバナーを使用したい場合は、手順について「[CPUfreq の設定](#)」を参照してください。
- ※ **--affected-cpus** — 周波数調整ソフトウェアを必要とする CPU を一覧表示します。
- ※ **--policy** — 現在の CPUfreq ポリシーの範囲 (KHz 単位) と現在アクティブなガバナーを表示します。
- ※ **--hwlimits** — CPU 使用できる周波数を KHz 単位で一覧表示します。

cpupower frequency-set では、以下のオプションを使用することができます。

- ※ **--min <freq>** と **--max <freq>** — CPU の *ポリシーの限界* を KHz 単位で設定します。



重要

ポリシーの限界を設定する場合は、**--max** を先に設定してから **--min** を設定してください。

- ※ **--freq <freq>** — CPU に特定のクロック速度を KHz 単位で設定します。設定できる速度は CPU のポリシーの限界範囲内に限られます (**--min** と **--max**)。
- ※ **--governor <gov>** — 新しい CPUfreq ガバナーを設定します。



注記

`cpupowerutils` パッケージをインストールしていない場合は、CPUfreq の設定は `/sys/devices/system/cpu/[cpuid]/cpufreq/` 内にある調節可能値で確認することができます。たとえば、cpu0 の最小クロック速度を 360 KHz に設定する場合は次のコマンドを使用します。

```
echo 360000 >
/sys/devices/system/cpu/cpu0/cpufreq/scaling_min_freq
```

3.3. CPU モニター

cpupower には、アイドルとスリープ状態の統計値および周波数情報を提供しプロセッサのトポロジーに関してレポートを行なう各種のモニター機能が備わっています。プロセッサ固有のモニターもあれば、あらゆるプロセッサに互換性があるものもあります。各モニターの測定対象や互換性のあるシステムなどに関する詳細については、**cpupower-monitor** の man ページをご覧ください。

次のオプションは、**cpupower monitor** コマンドに付けて使用します。

- ※ **-l** — システムで使用できる全モニターを一覧表示します。
- ※ **-m <monitor1>, <monitor2>** — 特定のモニターを表示します。識別子については **-l** を実行して確認します。
- ※ **command** — 特定コマンドに関する CPU の需要とアイドル統計値を表示します。

3.4. CPU 節電ポリシー

`cpupower` を使うと、プロセッサの節電ポリシーが調整できます。

次のオプションは `cpupower set` コマンドに付けて使用します。

`--perf-bias <0-15>`

対応している Intel プロセッサ上のソフトウェアがよりアクティブに、最適なパフォーマンスと節電とのバランスを確定できるようにします。このオプションは他の節電ポリシーを上書きするものではありません。割り当てる値は 0 から 15 の間で、0 が最適なパフォーマンスとなり、15 なら最適な電力効率となります。

デフォルトでは、このオプションはすべてのコアに適用されます。コア別に適用する場合は、`-cpu <cpu list>` オプションを追加します。

`--sched-mc <0|1|2>`

他の CPU パッケージが選ばれるまで、一つの CPU パッケージ内のコアに対するシステムプロセッサの電力使用を制限します。0 は制限なし、1 は最初は CPU パッケージをひとつだけ採用、2 は 1 に加えてタスクの復帰を処理する場合にセミアイドルの CPU パッケージを優先します。

`--sched-smt <0|1|2>`

他のコアが選ばれるまで、一つの CPU コアの複数スレッドに対するシステムプロセッサの電力使用を制限します。0 は制限なし、1 は最初は CPU パッケージをひとつだけ採用、2 は 1 に加えてタスクの復帰を処理する場合にセミアイドルの CPU パッケージを優先します。

3.5. サスペンドと復帰

システムがサスペンド状態になると、カーネルはドライバーを呼び出してその状態を保存し、それからドライバーをアンロードします。システムが復帰する時には、ドライバーを再読み込みし、デバイス群を再プログラムします。このタスクを遂行するドライバーにより、システムが正常に復帰できるかどうかが決まります。

この点では、ビデオドライバーが特に問題です。その理由は、ACPI (電力制御インタフェース : *Advanced Configuration and Power Interface*) 規格では、システムファームウェアがビデオハードウェアを再プログラムできる必要がないためです。そのため、ビデオドライバーがハードウェアを完全な未初期化の状態からプログラムできない限りは、システムは復帰できないことがあります。

Red Hat Enterprise Linux 7 では、新しいグラフィックチップセットをより強力的にサポートしています。これにより、サスペンドと復帰は以前より多くのプラットフォームで機能します。特に、NVIDIA チップセットに対するサポートは格別に向上しており、GeForce 8800 シリーズでは特に改善されています。

3.6. Active-State Power Management

Active-State Power Management (ASPM) は、接続するデバイスが使用中でない時に PCIe リンク用に電力状態を低く設定することで *Peripheral Component Interconnect Express* (PCI Express または PCIe) サブシステムで電力を節約します。ASPM はリンクの両端で電力状態を制御して、リンク末端のデバイスで電力が最大の場合でもリンク内で電力を節約します。

ASPM が有効な時には、異なる電力状態の間でのリンクを切り替えるために時間が必要なため、デバイスの遅延は大きくなります。ASPM には、電力状態を決定する以下の 3 つのポリシーがあります。

デフォルト

システムのファームウェア (例えば、BIOS) で指定されたデフォルトに従って、PCIe リンクの電

ソフトウェアの管理 (例えば、BIOS) で指定されたノードに基づいて、PCIe リンクの電力状態を設定します。これが ASPM のデフォルト状態です。

powersave

パフォーマンスの低下に関係なく、できる限り電力を節約するように ASPM を設定します。

performance

PCIe リンクが最大パフォーマンスで稼働できるように ASPM を無効にします。

ASPM のサポートは `pcie_aspm` カーネルパラメータで有効にしたり無効にしたりすることができます。`pcie_aspm=off` を使うと ASPM は無効になり、`pcie_aspm=force` にすると ASPM は有効になります。ASPM に対応していないデバイス上でも使用できます。

ASPM のポリシーは `/sys/module/pcie_aspm/parameters/policy` で設定しますが、`pcie_aspm.policy` カーネルパラメータを使って起動時に指定することも可能です。例えば、`pcie_aspm.policy=performance` を使用すると ASPM の performance ポリシーに設定されます。



警告

`pcie_aspm=force` を設定すると、ASPM をサポートしていないハードウェアでは、システムが反応しなくなる恐れがあります。`pcie_aspm=force` を設定する前に、システム上のすべての PCIe ハードウェアが ASPM をサポートすることを確認してください。

3.7. Aggressive Link Power Management

Aggressive Link Power Management (ALPM) は、アイドル時 (I/O が存在しない時) にディスクへの SATA リンクを低電力に設定することにより、ディスクの節電を促進する節電技術です。I/O リクエストがそのリンクにキューされると、ALPM は SATA リンクを自動的にアクティブな電力状態に設定し直します。

ALPM で導入された節電には、ディスクの遅延が伴います。そのため、アイドル状態の I/O 時間が長くなると思われる場合にのみ ALPM を使用するべきです。

ALPM は、Advanced Host Controller Interface (AHCI) を使用する SATA コントローラ上でのみ利用できます。AHCI の詳細情報については、<http://www.intel.com/technology/serialata/ahci.htm> を参照してください。

利用可能時には、ALPM はデフォルトで有効になっています。ALPM には以下の 3 つのモードがあります。

min_power

このモードは、ディスクに I/O がいない時に最小電力状態 (SLUMBER) へのリンクを設定します。このモードはアイドル時間が長くなると思われる場合に役立ちます。

medium_power

このモードは、ディスク上に I/O がいない時に 2 番目に電力が低い状態へのリンクを設定します。このモードでは、パフォーマンスへの影響ができるだけ少なくなるように、リンクの電力状態を切り替えることができます (例えば、一時的に I/O が多くなりまたアイドルになる間)。

`medium_power` モードでは、負荷に応じてリンクが PARTIAL と電力が最大の ACTIVE 状態の間で切り替え可能になります。PARTIAL から SLUMBER に、そして PARTIAL に戻るリンクを直接切り替えることはできません。このような場合、どの電力状態も最初に ACTIVE 状態を経由せずに、他に切り替わることはできません。

max_performance

ALPM は無効です。ディスクに I/O が無い時は、リンクは低電力状態になりません。

ご使用の SATA ホストアダプタが実際に ALPM に対応しているかどうか確認するには、`/sys/class/scsi_host/host*/link_power_management_policy` ファイルが存在するかどうか確認します。設定を変更するには、このセクションに記載してある値をこのファイルに書き込むか、あるいはファイルを表示して現在の設定を確認します。



重要

ALPM を `min_power`、または `medium_power` に設定すると、自動的に「ホットプラグ」機能を無効にします。

3.8. Relatime ドライブアクセス最適化

POSIX 基準では、各ファイルが最後にアクセスされた時間を記録するファイルシステムのメタデータがオペレーティングシステムによって維持されていなければなりません。このタイムスタンプは **atime** と呼ばれ、これを維持するにはストレージに常時書き込みをする動作が必要になります。これらの書き込みにより、ストレージデバイスとそのリンクに常に電源が投入され、ビジー状態になります。**atime** データを使用するアプリケーションは少ないため、このストレージデバイスの動作が電力を浪費していることとなります。重要なことは、ストレージへの書き込みは、ファイルがストレージからではなくキャッシュから読み込まれた場合でも発生する点です。これまで、Linux カーネルでは **mount** 用の **noatime** オプションに対応してきたため、このオプションでマウントされたファイルシステムには **atime** データを書き込んでいませんでした。しかし、単に **atime** データを使用しないことにも問題があります。一部のアプリケーションは **atime** データに依存しているため、これが利用できないと機能しないためです。

Red Hat Enterprise Linux 7 で使用しているカーネルは、代替となる **relatime** に対応しています。**Relatime** では **atime** データを維持しますが、ファイルがアクセスされる度の書き込み動作はしません。このオプションを有効にすると、ファイルが変更された、つまり **atime** が更新された (**mtime**) 場合、またはファイルが最後にアクセスされてから一定以上の時間 (デフォルトでは 1 日) が経過している場合に限り、**atime** データがディスクに書き込まれます。

デフォルトでは、**relatime** が有効な状態ですべてのファイルシステムがマウントされるようになります。特定のファイルシステムに対してこのオプションを無効にしたい場合には、そのファイルシステムをマウントする際に **norelatime** オプションを使用します。

3.9. パワーキャッピング (Power Capping)

Red Hat Enterprise Linux 7 では、HP の *Dynamic Power Capping* (DPC) や Intel Node Manager (NM) テクノロジーなど、最近のハードウェアに見られるパワーキャッピング (電力制限) 機能を利用しています。パワーキャッピングにより、管理者はサーバーによる電力消費の上限を設定できるだけでなく、より効率的にデータセンターを計画できます。その理由は、既存の電力供給装置に過負荷をかけるリスクが大幅に減少するためです。また、管理者はさらに多くのサーバー群を同じ物理フットプリント (physical footprint) に配置でき、サーバーの電力消費が制限されると、確実に高負荷時に電力需要が利用可能な電力を超えないようになります。

HP Dynamic Power Capping

Dynamic Power Capping は、選ばれた ProLiant と BladeSystem のサーバーで利用できる機能であり、システム管理者が 1 つのサーバー、あるいはサーバーのグループの電力消費量を制限できるようにします。

キャップとは、現時点の作業負荷に関係なく、サーバーが超過しない確実な上限のことです。キャップには、サーバーがその消費電力の上限に到達するまでは何の効果もありません。到達した時点で、管理プロセッサは CPU P 状態 と クロックスロットル (clock throttling) を調節して消費電力を制限します。

Dynamic Power Capping は、オペレーティングシステムから独立して CPU の動作を個別に修正しますが、HP の *integrated Lights-Out 2 (iLO2)* ファームウェアにより、オペレーティングシステムは管理プロセッサにアクセスでき、その結果ユーザスペースのアプリケーションは管理プロセッサにクエリできます。Red Hat Enterprise Linux 7 で使用されているカーネルには HP iLO と iLO2 のファームウェア用のドライバーが含まれており、プログラムが `/dev/hpilo/dXccbW` で管理プロセッサにクエリできるようにします。カーネルには、パワーキャッピング機能をサポートするための `hwmon sysfs` インターフェースの拡張と、`sysfs` インターフェースを使用する ACPI 4.0 パワーメーター用の `hwmon` ドライバーが含まれています。これらの機能が一緒になって、オペレーティングシステムとユーザスペースのツールがパワーキャップ用に設定された値とシステムの現在の電力消費量を読み込めるようになります。

HP Dynamic Power Capping についての詳細情報

は、<http://h20000.www2.hp.com/bc/docs/support/SupportManual/c01549455/c01549455.pdf> にある『HP Power Capping and HP Dynamic Power Capping for ProLiant Servers』を参照してください。

Intel Node Manager

Intel Node Manager は、CPU パフォーマンス、ひいては電力消費量を制限するためにプロセッサの P 状態と T 状態を使用して、システムにパワーキャップをかけます。電源管理ポリシーを設定することにより、管理者は、例えば夜間や週末などのシステムの負荷が低い時に電力消費が低くなるよう設定することができます。

Intel Node Manager は、標準の *電力制御インターフェース (Advanced Configuration and Power Interface)* を通じて、OSPM *オペレーティングシステム向け構成および電力管理 (Operating System-directed configuration and Power Management)* を使用することで CPU のパフォーマンスを調整します。Intel Node Manager が OSPM ドライバーに T 状態への変更を通知すると、そのドライバーは P 状態に対応する変更を加えます。同様に Intel Node Manager が OSPM ドライバーに P 状態への変更を通知すると、ドライバーはそれに応じて T 状態を変更します。こうした変更は自動的に発生し、オペレーティングシステムの介入を必要としません。管理者は *Intel Data Center Manager (DCM)* ソフトウェアを使用して Intel Node Manager の設定と監視を行います。

Intel Node Manager についての詳細情報は、<http://communities.intel.com/docs/DOC-4766> にある『Node Manager — A Dynamic Approach To Managing Power In The Data Center』を参照してください。

3.10. 拡張グラフィックス電力管理

Red Hat Enterprise Linux 7 は不必要な消費が発生するソースを取り除くことにより、グラフィックスおよびディスプレイデバイスの節電を行います。

LVDS 再クロック

LVDS *低電圧差動信号 (Low-voltage differential signalling)* とは、電子信号を銅線上で伝えるシステムです。このシステムが応用されている重要な例の 1 つは、ピクセル情報をノート PC の *液晶ディスプレイ (LCD)* 画面に送信することです。すべてのディスプレイには *リフレッシュレート* があります。これはディスプレイがグラフィックコントローラから新しいデータを受け取り、画像を画面に再表示する頻度です。通常、画面は毎秒 60 回新しいデータを受信します (60 Hz の周波数)。画面とグラフィックコントローラが LVDS でリンクされている時は、LVDS システムはリフレッシュのたびに電力を使用します。アイドル状態の時、多くの LCD 画面のリフレッシュレートは、目立った変化なく 30 Hz まで低下することがあります (リフレッシュレートが低下すると特有のフリッカーが起こる *ブラウン管 (CRT) モニター* とは異なります)。Red Hat Enterprise Linux 7 のカーネルに組み込まれている Intel グラフィックスアダプタ用のドライバーは、自動的にこの *ダウニングクロック (downclocking)* を実行し、画面がアイドル状態の時には約 0.5 W の節電をします。

メモリのセルフリフレッシュの有効化

SDRAM *Synchronous dynamic random access memory* — これは、グラフィックスアダプタのビデオメモリに使用されます。毎秒何千回もリチャージされるため、個々のメモリセルは保管されているデータを保持します。データはメモリの内外へと移動するためそのデータを管理するその主要機能の他に、メモリコントローラには通常これらのリフレッシュサイクルを開始する役割があります。一方、SDRAMには低電力のセルフリフレッシュモードもあります。このモードでは、メモリは内部タイマーを使用して、そのリフレッシュサイクルを生成します。これにより、現在メモリに保存されているデータを危険にさらすことなく、システムはメモリコントローラをシャットダウンできます。Red Hat Enterprise Linux 7で使用されているカーネルは、アイドル状態の時に Intel グラフィックスアダプタのメモリにセルフリフレッシュをさせることができます。これにより約 0.8 W の節電ができます。

GPU クロックの低減

標準的なグラフィカルプロセッシングユニット (GPU) には、その内部回路の各種パーツを制御する内部クロックが含まれています。Red Hat Enterprise Linux 7で使用されているカーネルは、Intel および ATI の GPU 内の内部クロックの一部の周波数を低くすることができます。GPU コンポーネントが所定時間内に実行するサイクル数を低減すると、それらが実行する必要がなかったサイクルで消費されていたであろう電気を節減します。GPU がアイドル状態の時には、カーネルは自動的にそうしたクロックの速度を遅くし、GPU の活動が増加すると速めます。GPU のクロックサイクルを低下させることで、最大で約 5 W の節電ができます。

GPU の電源オフ

Red Hat Enterprise Linux 7 の Intel と ATI グラフィックスドライバーは、アダプタにモニターが接続されていない時を検出できるため、GPU を完全にシャットダウンすることができます。この機能は、常時モニターを接続していないサーバーで特に重要です。

3.11. RFKill

多くのコンピュータシステムには、Wi-Fi、Bluetooth、および 3 G デバイスを含む無線送信器が搭載されています。これらのデバイスは電力を消費し、使用していない時には無駄になります。

RFKill は、コンピュータシステムの無線送信器が、クエリ、アクティベート、非アクティブ化されるインターフェースを提供する Linux カーネルのサブシステムです。無線送信器が非アクティブ化されると、それらはソフトウェアが再アクティベートできる状態 (ソフトブロック) に置かれるか、またはソフトウェアが再アクティベートできない状態 (ハードブロック) に置かれます。

RFKill コアは、サブシステムにアプリケーションプログラミングインターフェース (API) を提供します。RFKill をサポートするように設計されているカーネルドライバーは、この API を使用してカーネルに登録します。また、デバイスを有効および無効にする方法を含んでいます。さらに RFKill コアは、ユーザーアプリケーションが解釈できる通知と、ユーザーアプリケーションが送信器の状態をクエリする方法を提供します。

RFKill インターフェースは `/dev/rfkill` にありますが、システムのすべての無線送信器の現在の状態が含まれています。各デバイスの現在の RFKill の状態は、`sysfs` に登録されています。また、RFKill は RFKill 対応のデバイス内の状態の変化について `uevents` を発行します。

`Rfkill` は、システム上の RFKill 対応のデバイスをクエリ、変更できるコマンドラインツールです。このツールを取得するには、`rfkill` パッケージをインストールしてください。

コマンド `rfkill list` を使用すると、デバイスの一覧が取得できます。それぞれのデバイスにはそれに関連した `0` から始まるインデックス番号があります。このインデックス番号を使用して `rfkill` に対してデバイスのブロックとブロック解除を指示します。例を示します。

```
rfkill block 0
```

上記は、システムの最初の RFKill 対応デバイスをブロックします。

また、**rfkill** を使用してデバイスの特定のカテゴリ、またはすべての RFKill 対応のデバイスもブロックできます。例を示します。

```
rfkill block wifi
```

システムのすべての Wi-Fi デバイスをブロックします。すべての RFKill 対応デバイスをブロックするには、以下を実行します。

```
rfkill block all
```

デバイスをブロック解除するには、**rfkill block** の代わりに **rfkill unblock** を実行します。**rfkill** がブロックできるデバイスカテゴリの全一覧を取得するには、**rfkill help** を実行してください。

第4章 使用例

この章では、2 種類のユースケースを使ってこのガイドで説明している分析と設定方法を説明しています。最初は、標準的なサーバーを例にとり、その次は標準的なノート PC を例にして考えてみます。

4.1. 例 — サーバー

今日の一般的な標準サーバーは、Red Hat Enterprise Linux 7 でサポートされている必要なハードウェアの機能がすべて搭載されています。最初に考慮すべきなのは、サーバーの主要な使用目的となる作業負荷の種類です。この情報に基づき、節電のためにどのコンポーネントを最適化するか決定できます。

サーバーのタイプに関係なく、グラフィックス性能は一般的には必要ありません。そのため、GPU 節電はオンのままで結構です。

ウェブサーバー

ウェブサーバーにはネットワークとディスク I/O が必要です。外部の接続スピードによっては、100 Mbit/s で十分かも知れません。マシンがほとんど静的なページを使用する場合は、CPU のパフォーマンスはあまり重要ではないでしょう。以下のような電力管理の選択肢があります。

- ※ **tuned** にはディスクまたはネットワークのプラグインなし。
- ※ ALPM をオンにする
- ※ **ondemand** ガバナーをオンにする
- ※ ネットワークカードは 100 Mbit/s に制限する

計算サーバー

計算サーバーには主に CPU が必要です。以下のような電力管理の選択肢があります。

- ※ ジョブとデータストレージが発生する場所に応じて、**tuned** のディスク、又はネットワークプラグイン。またはバッチモードシステムには、完全にアクティブな **tuned**。
- ※ 使用量によっては、**performance** ガバナー。

メールサーバー

メールサーバーには、多くの場合ディスク I/O と CPU が必要です。以下のような電力管理の選択肢があります。

- ※ **ondemand** ガバナーはオン。CPU パフォーマンスの最後の数パーセントは重要でないためです。
- ※ **tuned** にはディスクまたはネットワークのプラグインなし。
- ※ メールは内部で発生することが多く、1 Gbit/秒 か 10 Gbit/秒 のリンクから利用できるためネットワークスピードは制限しません。

ファイルサーバー

ファイルサーバーの要件はメールサーバーの要件に似ています。しかし使用するプロトコル次第では、さらなる CPU パフォーマンスが必要になる可能性があります。一般的に Samba ベースのサーバーは、NFS よりも CPU を要求して、NFS は一般的に iSCSI よりも CPU を要求します。それでも、**ondemand** ガバナーを使用できるはずで

ディレクトリーサーバー

ディレクトリーサーバーのディスク I/O の要件は、一般的に低いものです。十分な RAM がある場合は特にそうです。ネットワーク遅延は重要ですが、ネットワーク I/O はそれほどでもありません。リンクの速度が遅い遅延のネットワークのチューニングを考えられるかも知れませんが、これを特定のネットワークに注意深くテストするようにしてください。

4.2. 例 — ノート PC

電力管理と節電が実際に効果をもたらすもうひとつの非常に一般的な対象は、ノート PC です。ノート PC はもともとワークステーションやサーバーよりも大幅に少ないエネルギーを使用するように設計されているため、絶対的な節電ができる可能性は他のマシンよりも低くなります。ただし、バッテリーモードでは、どんな節電でもノートパソコンのバッテリー寿命を数分でも延長するのに役立ちます。このセクションでは、ノート PC のバッテリーモードにフォーカスしていますが、もちろん AC 電源での使用でもこうしたチューニングの一部、またはすべてを活用することができます。

1つのコンポーネントの節電は、通常ワークステーションよりもノートパソコンで相対的に大きな効果をもたらします。例えば、100 Mbits/秒 で実行している 1 Gbit/秒 ネットワークインターフェースはおよそ 3–4 ワット節約します。約 400 ワットの合計消費電力を持つ標準的なサーバーには、この節約はおよそ 1% です。約 40 ワットの合計消費電力を持つノートパソコンでは、この1つのコンポーネントの節電は合計でおよそ 10% になります。

標準的なノート PC での特定の節電最適化としては以下のものがあります。

- ※ システムの BIOS を使用しないすべてのハードウェアを無効にするように設定します。例えば、パラレルポートまたはシリアルポート、カードリーダー、Web カメラ、WiFi および Bluetooth などが可能です。
- ※ スクリーンを見るために最高輝度が必要ない暗めの場所では、ディスプレイ輝度を低くします。そのためには、GNOME デスクトップでは、システム+設定 → 電力管理 と進みます。KDE デスクトップでは、アプリケーション起動キックオフ (Kickoff Application Launcher) + コンピュータ+システム設定+高度な設定 → 電力管理 と進みます。または、コマンドラインで `gnome-power-manager` か、`xbacklight` を実行するか、ノート PC でファンクションキーを使用します。

また、(代わりに) 各種システム設定を微調整することもできます。

- ※ `ondemand` ガバナーを使用します (Red Hat Enterprise Linux 7 ではデフォルトで有効です)。
- ※ AC97 オーディオ節電機能を有効にします (Red Hat Enterprise Linux 7 ではデフォルトで有効です)。

```
echo Y > /sys/module/snd_ac97_codec/parameters/power_save
```

- ※ USB 自動サスペンドを有効にします。

```
for i in /sys/bus/usb/devices/*/power/autosuspend; do echo 1 > $i; done
```

USB 自動サスペンドはすべての USB デバイスで正常に機能するわけではありません。

- ※ `relatime` を使用してファイルシステムをマウントします (Red Hat Enterprise Linux 7 ではデフォルトです)。

```
mount -o remount,relatime mountpoint
```

- ※ 画面の輝度を 50 かそれ以下に下げます。例えば以下のとおりです。

```
xbacklight -set 50
```

- ※ スクリーンのアイドル状態に DPMS をアクティベートします。

```
xset +dpms; xset dpms 0 0 300
```

- ※ Wi-Fi を非アクティブ化します。

```
echo 1 > /sys/bus/pci/devices/*/rf_kill
```


付録A 開発者へのヒント

すべての優れたプログラミング教本では、メモリ割り当ての問題と特定機能のパフォーマンスを網羅しています。ご自身のソフトウェアを開発するにあたっては、ソフトウェアが実行するシステムで電力消費が増加する可能性があることに注意してください。こうした配慮は、コードのすべてのラインに影響するものではありませんが、頻繁にパフォーマンスのボトルネックになる領域ではコードを最適化することができます。

問題なることが多い手法は以下のとおりです。

- ✦ スレッドを使用する。
- ✦ 不必要な CPU のウェイクアップ、およびウェイクアップを効率良く使用しない状態。ウェイクアップする必要がある場合は、すべてを一度にできるだけ迅速に実行します (すぐにアイドル状態になるように迅速にすべてを実行します)。
- ✦ `[f]sync()` を不必要に使用する。
- ✦ 不必要なアクティブポーリング、または短い通常のタイムアウトを使用する (代わりにイベントに反応する)。
- ✦ ウェイクアップを効率的に使用していない。
- ✦ 低効率なディスクアクセス。頻繁なディスクアクセスを回避するために大きなバッファを使用してください。一度に大きなブロックを書き込みます。
- ✦ 低効率のタイマーを使用する。可能な限りアプリケーション群 (またはシステム群) にタイマーをグループ化します。
- ✦ 過度の I/O、電力消費、またはメモリ使用 (メモリリークを含む)。
- ✦ 不必要に計算を実行する。

以下のセクションでは、これらの領域をさらに詳しく検証していきます。

A.1. スレッドの使用

スレッドを使用するとアプリケーションのパフォーマンスが改善し、より速くなると思われていますが、これはすべてのケースで当てはまるわけではありません。

Python

Python は Global Lock Interpreter ^[1] を使用するため、スレッドは大規模な I/O 運用でのみ効果的です。Unladen-swallow ^[2] は、コードを最適化できる可能性がある Python の高速実装です。

Perl

Perl のスレッドは、もともとはフォークがないシステム (32-bit Windows オペレーティングシステムのシステムなど) で実行するアプリケーション用に開発されました。Perl のスレッドでは、データはすべての単独スレッド用にコピー (コピーオンライト: Copy On Write) されます。ユーザーはデータ共有のレベルを定義できるため、データはデフォルトでは共有されません。データを共有するには、`threads::shared` モジュールを含める必要があります。しかし、データはコピー (コピーオンライト) されるだけでなく、モジュールはデータのタイ変数も作成します。これでさらに時間がかかり、遅くなります。 ^[3]

C

C のスレッドは同じメモリを共有します。各スレッドはそれ自身のスタックを持ち、カーネルは新しい

ファイル記述子を作成したり、新しいメモリスペースの割り当てをする必要がありません。Cはより多くのスレッドにより多くのCPUのサポートを実際に活用できます。そのため、スレッドのパフォーマンスを最大にするには、CかC++などの低水準言語を使用すべきです。スクリプト言語を使用している場合は、Cバインディングを考慮してください。プロファイラを使用すると、コードの正しく実行していない部分を特定できます。 [4]

A.2. ウェイクアップ (Wake-ups)

多くのアプリケーションは、設定ファイルの変更を確認するためにスキャンします。多くの場合、スキャンは例えば毎分など、決まった間隔で実行されます。このスキャンが問題になる理由は、スキャンにより回転が停止しているディスクをウェイクアップさせるためです。最善策は、適切な間隔、適切な確認方法を見つけるか、**inotify** で変更を確認して、イベントに対応することです。**inotify** はファイルまたはディレクトリで様々な変更を確認できます。

例を示します。

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/time.h>
#include <sys/types.h>
#include <sys/inotify.h>
#include <unistd.h>

int main(int argc, char *argv[]) {
    int fd;
    int wd;
    int retval;
    struct timeval tv;

    fd = inotify_init();

    /* checking modification of a file - writing into */
    wd = inotify_add_watch(fd, "./myConfig", IN_MODIFY);
    if (wd < 0) {
        printf("inotify cannot be used\n");
        /* switch back to previous checking */
    }

    fd_set rfd;
    FD_ZERO(&rfd);
    FD_SET(fd, &rfd);
    tv.tv_sec = 5;
    tv.tv_usec = 0;
    retval = select(fd + 1, &rfd, NULL, NULL, &tv);
    if (retval == -1)
        perror("select()");
    else if (retval) {
        printf("file was modified\n");
    }
    else
        printf("timeout\n");

    return EXIT_SUCCESS;
}
```

このアプローチの優れている点は、確認できる方法が多岐に渡っている点です。

主な制限は、1つのシステムでは利用できる監視の数が限られていることです。この数は `/proc/sys/fs/inotify/max_user_watches` から取得できます。この数値を変更することは可能ですが、推奨されません。さらに、`inotify` が失敗すると、コードは別の確認方法にフォールバックする必要がありますが、これは通常ソースコードの `#if #define` が多く発生することを意味しています。

`inotify` の詳細については、`inotify(7) man` ページを参照してください。

A.3. Fsync

`Fsync` は I/O 負荷の高い動作として知られていますが、実際にはそうでない場合もあります。

ユーザーが新しいページに移動するためのリンクをクリックする度、`Firefox` は `sqlite` ライブラリを呼び出していました。`sqlite` が `fsync` を呼び出すため、そのファイルシステム設定 (主に `data-ordered` モードの `ext3`) が原因で、何も起こらない場合は長い待ち時間が発生していました。同時に別のプロセスが大きなファイルをコピーしている場合には、最大で 30 秒もの時間がかかっていました。

ただし、`fsync` が全く使用されていない別のケースでは、`ext4` ファイルシステムへの切り替えで問題が発生していました。`Ext3` は `data-ordered` モードに設定され、数秒毎にメモリーがフラッシュされ、その内容がディスクに保存されていました。ただし、`ext4` の `laptop_mode` では、保存の間隔が長いため、システムの電源が予期せずオフになった場合にデータが消失する可能性がありました。現在、`ext4` にはバッチが適用されましたが、アプリケーションの設計を慎重に考慮し、適切に `fsync` を使用する必要があります。

設定ファイルからの読み込みと設定ファイルへの書き込みに関する簡単な例を使って、ファイルのバックアップが作成される流れと、データが消失してしまう流れを示します。

```
/* open and read configuration file e.g. ./myconfig */
fd = open("./myconfig", O_RDONLY);
read(fd, myconfig_buf, sizeof(myconfig_buf));
close(fd);
...
fd = open("./myconfig", O_WRONLY | O_TRUNC | O_CREAT, S_IRUSR |
S_IWUSR);
write(fd, myconfig_buf, sizeof(myconfig_buf));
close(fd);
```

より適切な例は以下のようになります。

```
/* open and read configuration file e.g. ./myconfig */
fd = open("./myconfig", O_RDONLY);
read(fd, myconfig_buf, sizeof(myconfig_buf));
close(fd);
...
fd = open("./myconfig.suffix", O_WRONLY | O_TRUNC | O_CREAT, S_IRUSR |
S_IWUSR);
write(fd, myconfig_buf, sizeof(myconfig_buf));
fsync(fd); /* paranoia - optional */
...
close(fd);
rename("./myconfig", "./myconfig~"); /* paranoia - optional */
rename("./myconfig.suffix", "./myconfig");
```

[1] <http://docs.python.org/c-api/init.html#thread-state-and-the-global-interpretor-lock>

[2] <http://code.google.com/p/unladen-swallow/>

[3] http://www.perlmonks.org/?node_id=288022

[4] <http://people.redhat.com/drepper/lt2009.pdf>

付録B 改訂履歴

改訂 2-0.1	Thu Sep 3 2015	Takuro Nagamoto
翻訳ファイルを XML ソースバージョン 2-0 と同期		
改訂 2-0	Wed 18 Feb 2015	Jacquelynn East
7.1 GA 向けバージョン		