

ASR5x00 MME オーバーロード保護機能

目次

[概要](#)

[MME 保護](#)

[ネットワークの過負荷保護：アタッチ率スロットリング](#)

[ネットワークの過負荷保護：ページングスロットリング](#)

[設定例](#)

[ネットワークの過負荷保護：DDN スロットリング \(サービング ゲートウェイ機能、MME 保護 \)](#)

[ネットワークの過負荷保護：EGTP パス障害のスロットリング](#)

[設定例](#)

[拡張輻輳制御](#)

[輻輳状態のしきい値](#)

[しきい値と許容レベル](#)

[サービス コントロール CPU しきい値](#)

[システム CPU しきい値](#)

[システム メモリしきい値](#)

[関連情報](#)

概要

このドキュメントでは、Cisco アグリゲーション サービス ルータ (ASR) 5000 シリーズで使用できる、モビリティ マネージメント エンティティ (MME) のさまざまな過負荷保護方式と機能に焦点を当てます。ASR 5000 シリーズはユーザによる制御を実現するさまざまな方法を備えています。この記事は、そうした機能および関連する CLI コマンドについて説明します。

MME 保護

ネットワークの過負荷保護：アタッチ率スロットリング

アタッチ率スロットリングはホーム サブスクライバ サーバ (HSS)、ポリシー/課金ルール機能 (PCRF)、オンライン課金サーバ (OCS) などの隣接ネットワーク要素と、imsimgr や sessmgr などの内部 MME リソースを保護します。アタッチ率スロットリングは、Attach および Inter-MME/サービング GPRS サポート ノード (SGSN) トラッキング エリア更新 (TAU) など、imsimgr に到着した新しいコールを処理します。

この図は、コールとスロットリング キューのメッセージ フローを示しています。

MME (imsimgr と sessmgr 前方) を保護するには、スロットリング率、キュー待機時間、キューサイズ時間を定義する必要があります。スロットリング率は MME コール モデルによって異なります。これは MME キャパシティがコール モデルによって異なるためです。

MME の場合、スロットリング率の計算は比較的単純で、ネットワーク内の標準的な 1 秒あたりのコール イベント数 (CEPS) に許容度を加えます。また、HSS の保護も必要であれば、HSS データベース容量も考慮する必要があります。

例

処理が多い時間帯に、MME は最大で毎秒 170 ~ 200 件のコールを処理します (Attach + Inter TAU)。あるサイトで障害が発生すると、毎秒 350 ~ 370 件のコールが 1 つの MME に到達する可能性があります。このようなコール率の場合、MME 使用率は 80% 近くまで上がり、MME ボックス内の過剰なシグナリング負荷を回避するためには、毎秒 400 件のコールにスロットリング率を制限するのが最適なレベルです。

キュー待機時間はデフォルトで 5 秒です。これは「お客様」にとって最適です。キューサイズはデフォルトで 2500 秒です。これは「お客様」にとって最適です。

設定コマンドは、次のようになります。

```
asr5k(config)#network-overload-protection mme-new-connections-per-second
new_connections action attach { drop | reject-with-emm-cause
{ congestion | network-failure | no-suitable-cell-in-tracking-area}
tau { drop | reject-with-emm-cause { congestion | network-failure
| no-suitable-cells-in-tracking-area | no-sec-ctxt-in-nw} fwd-reloc
{ drop | reject} }{wait-time <wait-time>} {queue-size <queue-size>}
```

new_connections

毎秒受け入れる新しい MME 接続の数を定義します。これは、50 ~ 5000 の整数にする必要があります。デフォルトは 500 です。

action

ペーシング キューがいっぱいになったときにとるアクションを定義します。MME で新しい接続が受信されるたびに、それらはペーシング キューにキューイングされ、imsimgr は設定済みレートでキューからのメッセージを処理します。(高い着信レートのために) キューがオーバーフローすると、設定された「アクション」に基づいてパケットがドロップされるか、または拒否されます。

queue-size

パケットのバッファリングに使用されるペーシング キューの最大サイズを定義します。これは、250 ~ 25000 の整数にする必要があります。デフォルトは 2500 です。

設定例

```
network-overload-protection mme-new-connections-per-second 400 action attach
reject-with-emm-cause no-suitable-cell-in-tracking-area tau
reject-with-emm-cause no-suitable-cell-in-tracking-area fwd-reloc drop
```

ここで 1 秒あたりのコール率が 400 に設定され、アクションは原因 #15 を伴うインテリジェント拒否であるため、ユーザ機器 (UE) は別の無線アクセス技術 (RAT) に再接続されます。待機時間はデフォルト (5 秒) に設定され、キューサイズは 2500 です。

注: EMM 原因 #15 (no-suitable-cell-in-tracking-area) を伴うアクション「拒否」が適切な

理由は、#15 で拒否されたコールのほとんどは MME に再到達することがなく、別の RAT 層 (3G、2G) に送られるためです。サービング無線ネットワークサブシステムリロケーション (SRNS) に関するアクション「ドロップ」は将来使用されるもので、拒否後に MME にすぐに再アタッチするのを防ぎます。

ネットワークの過負荷保護：ページングスロットリング

ページングスロットリングは内部 MME リソース (mmemgr) を必要に応じて eNodeB/radio リソースとして保護します。このレート制限しきい値は、特定の ASR 5000 シャーシで MME に関連付けられるすべての eNodeB に適用されます。eNodeB への S1 ページング要求は、このしきい値でレート制限されます。このしきい値を超える eNodeB への S1 ページング要求はドロップされます。

MME の場合、スロットリング率の計算は比較的単純で、ネットワーク内の標準的な出力ページング率に許容度を加えます。(これは純粋に、設計チームの決定に基づきます)。

例

処理が多い時間帯に、各 MME は毎秒最大 35000 件のページングメッセージを処理します。あるサイトで障害が発生すると、毎秒最大 70000 ページが 1 つの MME から送信される可能性があります。このようなページング率の場合、MME の使用率 (mmemgr) は 80% 近くまで上がり、mmemgr を介した過剰な S1 シグナリングを回避するために、毎秒 70000 ~ 80000 ページにスロットリング率を制限するのが最適なレベルです。

ただし、このレート (率) は平均的な eNodeB によって制限されます。eNodeB ごとの平均レート (6500 eNodeB の場合) は、毎秒 10 ページです。ただし、トラッキングエリア (TA) はサブスクライバ数で同じではなく、さまざまな TA/メンバーの eNodeB が異なったページングでロードされます。TA ごとの平均サブスクライバ数に対して TA のサイズが 2 倍異なる場合は、eNodeB ごとのレートが 20 になります。TA ごとの平均サブスクライバ数に対して TA のサイズが 20 倍違う場合は、eNodeB ごとのレートが 200 になります。つまり、(サブスクライバ数の) TA が均等にロードされる場合に、この機能が最も効率的になります。

並行して行うべき別のアクションは、インテリジェント ページングをアクティブにすることです。『ASR 5000 MME Administration Guide』の「TAI mgmt db および LTE ページング」の項を参照してください。

設定コマンドは次のようになります。

```
asr5000(config)# network-overload-protection mme-tx-msg-rate-control enb s1-paging  
<rate in messages per second>
```

- network-overload-protection は、ネットワーク過負荷保護を示します
- mme-tx-msg-rate-control enb は平均的な eNodeB ごとの MME メッセージ レート制御を示します
- s1-paging は S1 ページング用のメッセージ レート制御を示します
- <rate> は eNodeB ごとの 1 秒あたりのメッセージのレートしきい値を指定します (1 ~ 65535)。

設定例

```
asr5000(config)# network-overload-protection mme-tx-msg-rate-control enb s1-paging
<rate in messages per second>
```

注：

- レート制限は、減少させる方向でさらに調整すべき項目です。TA を通るサブスクリバ数 (ページング数) に基づいて調整します (TA レベルの統計情報が必要です) 。
- この機能は (TA ごとのサブスクリバ/ページング数の点で) TA が均等にロードされる場合に、最も効率的になります。

ネットワークの過負荷保護：DDN スロットリング (サービング ゲートウェイ機能、MME 保護)

ダウンリンク データ通知 (DDN) スロットリングは、サービング ゲートウェイ (SGW) 側から MME への DDN 要求のレートを制御する機能です。これは DDN (つまり入力ページング要求) のサージから mmemgr や sessmgr などの MME リソースを保護します。

この機能には 2 つの部分があります。1 つは Rel-10 準拠 MME 用で、もう 1 つは Rel-10 非準拠 MME 用です：

- Rel-10 準拠 MME の場合、この機能を有効にするには、SGW サービスの DDN スロットリング割り当ておよび保持プライオリティ (ARP) 上限を設定します。
- Rel-10 非準拠 MME の場合、SGW サービスでの ARP 上限に加えて他のパラメータをいくつか設定する必要があります (スロットリング要因、スロットリング時間、安定化時間、ポーリング間隔など) 。

SGW でこの機能を有効にすると、DDN 要求で ARP 上限が MME に送信されます。応答では、MME がスロットリング遅延単位、スロットリング遅延値、およびスロットリング要因を送信します。遅延値と遅延単位の組み合わせでスロットリング時間が計算されます。これらの値を受信すると、SGW はスロットリング タイマーが切れるまで、特定の ARP の DDN 要求をドロップします。

ローカル設定を使用する非 Rel-10 準拠 MME では、SGW は特定の ARP 上限を使って DDN 要求をスロットリングします。

Cisco ASR5x00 MME リリース 16 および 17 は自動 DDN スロットリングをサポートしていません。そのため、DDN スロットリングに関しては非 Rel 10 準拠として機能します。

注: DDN スロットリングは、出力側 (S1) ではなく入力側 (S11) で MME ページング スロットリングに加えてさらにきめ細かなスロットリングを提供します。ページング スロットリングを設定済みの場合、DDN スロットリングの実装は必須ではありませんが、これを実装するとより早い段階で過負荷を検出して対処できます。

技術仕様 (TS) 23.401、MME に関する参照項目：

DDN 要求のスロットリング

まれな状況下 (オペレータの設定したしきい値を MME ロードが超えた場合など) では、SGW が生成するシグナリング ロードを MME が制限することがあります (そのように設定されている場合) 。

MME はアイドル モードの UE に関する優先度が低いトラフィックの DDN リクエストを拒否したり、さらに MME をオフロードしたりできます。MME は、スロットリング係数に従い、DDN Ack メッセージで指定されたスロットリング遅延に基づき、アイドル モード UE に関して受信されたダウンリンク低優先度トラフィック用に送信する DDN 要求の数を選択的に減らすように SGW に要求できます。

SGW は、ベアラの ARP 優先度レベルとオペレータ ポリシー (優先/非優先トラフィックと見なすための ARP 優先度レベルに関する SGW でのオペレータ設定) に基づき、そのベアラが低優先度トラフィック用かどうかを判別します。MME は、SGW から受信した ARP 優先度レベルとオペレータ ポリシーに基づいて、低優先度トラフィック用の DDN 要求かどうかを判別します。

アイドル状態シグナリング低減 (ISR) が UE に関してアクティブでない場合、スロットリング遅延時に SGW は、その MME によって処理されるユーザ プレーン非接続と認識される (つまり SGW コンテキスト データにダウンリンク ユーザ プレーン トンネル端点識別子 (TEID) が示されない) UE に関するすべての低優先度ベアラで受信されたダウンリンク パケットを、スロットリング係数に基づいてドロップし、スロットリング対象外のベアラだけに関して DDN メッセージを MME に送信します。

スロットリング遅延時に UE に関する ISR がアクティブな場合、SGW は MME に DDN を送信せず、SGSN にのみ DDN を送信します。MME および SGSN の両方が負荷低減を要求した場合、SGW は、ユーザ プレーン非接続と認識される (つまり SGW コンテキスト データにダウンリンク ユーザ プレーン TEID が示されない) UE に関するすべての低優先度ベアラで受信されたダウンリンク パケットを、スロットリング係数に基づいてドロップします。

スロットリング遅延が期限切れになると、SGW は通常のコアネットワークを再開します。最後に受信したスロットリング係数値およびスロットリング遅延は、その MME から受信した以前のすべての値に代わって使用されるようになります。スロットリング遅延の受信により、その MME に関連付けられた SGW タイマーが再び開始します。

SGW の場合、MME と比べてスロットリング率の計算が比較的単純です。許容される最大入力ペーシング率 (MME ボックスごとに毎秒 1100 件のメッセージ) を考慮します。

コンフィギュレーション コマンドは次のとおりです。

```
asr5000(config)# network-overload-protection mme-tx-msg-rate-control enb s1-paging
<rate in messages per second>
```

throttle arp-watermark arp_value

ARP 上限が設定されている場合、MME/SGSN が DDN ACK メッセージでスロットリング係数と遅延を送信すると、設定値よりも大きい ARP 値を持つすべての DDN は、指定された遅延とスロットリング係数に基づいてスロットリングされます。

arp_value は、1 から 15 までの整数です。

rate-limit limit

レート制限を設定します (これと、後続のトークンを使用することで、MME がリリース 10 MME でない場合にのみレート制限します)。

limit は、1 から 999999999 の整数です。

time-factor seconds

SGW がスロットリングの決定を行う時間の長さを設定します。

seconds は、1 から 300 の整数です。

throttle-factor percent

DDN スロットリング係数を設定します。DDN 急増が検出された際にドロップされる DDN の割合を入力します。

percent は、1 から 100 の整数です。

increment-factor percent

DDN スロットリングの増分係数を設定します。DDN スロットリングを増分するパーセンテージを入力します。

percent は、1 から 100 の整数です。

poll-interval seconds

DDN スロットリングのポーリング間隔を設定します。

seconds は、2 から 999999999 の整数です。

throttle-time-sec seconds

DDN スロットリングの時間を秒数で設定します。SGW で DDN がスロットリングされる期間を秒数で入力します。

seconds は、0 から 59 の整数です。

throttle-time-min minutes

DDN スロットリング時間を分単位で設定します。SGW で DDN がスロットリングされる期間を分単位で入力します。

minutes は、0 から 59 の整数です。

throttle-time-hour hour

DDN スロットリング時間を時間数で設定します。SGW で DDN がスロットリングされる期間を時間数で入力します。

hour は、0 から 310 の整数です。

stab-time-sec seconds

DDN スロットリングの安定化時間を秒数で設定します。その期間 (秒) 内にシステムが安定化したら、スロットリングを無効にする期間を秒数で入力します。

seconds は、0 から 59 の整数です。

stab-time-min minutes

DDN スロットリングの安定化時間を分単位で設定します。その期間 (分) 内にシステムが安定化したら、スロットリングを無効にする期間を分単位で入力します。

minutes は、0 から 59 の整数です。

stab-time-hour hour

DDN スロットリングの安定化時間を時間数で設定します。その期間 (時間) 内にシステムが安定化したら、スロットリングを無効にする期間を分単位で入力します。

hour は、0 から 310 の整数です。

設定例

```
ddn throttle arp-watermark 1 rate-limit RATE time-factor 120 throttle-factor 50
increment-factor 10 poll-interval 30 throttle-time-sec 0 throttle-time-min 1
throttle-time-hour 0 stab-time-sec 0 stab-time-min 2 stab-time-hour 0
```

- 許容される最大入力レートは 1100 ページ/秒です (DDN を含む)
- DDN 急増時の 1100 ページ/秒は、1100 DDN/秒に相当します
- MME サイトごとに 4 つの SGW を備えた地域の場合 > レート = 275 DDN/秒 (SGW あたり)

-) の最大レートが可能
- MME サイトごとに 3 つの SGW を備えた地域の場合 > レート = 366 DDN/秒 (1 SGW あたり) の最大レートが可能
- MME サイトごとに 2 つの SGW を備えた地域の場合 > レート = 550 DDN/秒 (1 SGW あたり) の最大レートが可能
- MME サイトごとに 1 つの SGW を備えた地域の場合 > レート = 1100 DDN/秒 (1 SGW あたり) の最大レートが可能

ネットワークの過負荷保護：EGTP パス障害のスロットリング

この機能は、IP バックボーンと IP バックホールでの送信障害時、およびサイド ネットワーク要素の障害/再起動時に、拡張 GPRS トンネリング プロトコル (EGTP) パス障害の急増から MME リソース (sessmgr、mmemgr) および 4G リソースを保護します。この機能は、検出された EGTP パス障害イベントを sessmgr ごとに制限することを可能にし、S1 ページング スロットリングに加えてさらにきめ細かなサブスライバ管理を定義します。アイドル状態のサブスライバと接続されたサブスライバの分離に基づいて、制限が設定されます。これは極めてネットワーク固有であるため、eUTRAN および UE ステータスに関連した調整が必要です。

例

サブスライバはおよそ 80:20 の割合、アイドル状態と接続状態に分離されます。最悪の場合、アイドル状態サブスライバに関する EGTP PF によりページングが急増し、それが原因で mmemgr の過負荷、つまりチェーン内の最も狭いボトルネックが発生する可能性があります。(アイドル状態サブスライバに関する)このような EGTP ページング係数 (PF) 急増により、まずページングが急増し、その急増が mmemgr ボトルネックにつながるため、まずこの状態に対して mmemgr を保護する必要があります。こうして IDLE の EGTP PF を予期しない出力ページング急増と見なすことができ、その許容値は最大で 1 秒あたり 1100 ページとなります。

- つまり、アイドル サブスライバに関する推奨されるスロットリング制限は、1 秒あたり 1000 メッセージです。
- 接続された sub の数はあります | 5 から 7 回 IDLE よりより少し。
- 接続されたサブスライバではページングの急増が発生しないため、接続されたサブスライバに関して安全に適用できる推奨値は 2000 メッセージ/秒です。

注: EGTP PF スロットリングは、出力側 (S1) ではなく入力側 (S11、Sv) で MME ページング スロットリングに加えてさらにきめ細かなスロットリングを提供します。ページング スロットリングを設定済みの場合、EGTP PF スロットリングの実装は必須ではありませんが、これを実装するとより早い段階で過負荷を検出して対処できます。

この設定は、インターフェイス タイプ 「interface-mme」である EGTP サービスに適用されます。

設定コマンドは次のようになります。

```
ddn throttle arp-watermark 1 rate-limit RATE time-factor 120 throttle-factor 50
increment-factor 10 poll-interval 30 throttle-time-sec 0 throttle-time-min 1
throttle-time-hour 0 stab-time-sec 0 stab-time-min 2 stab-time-hour 0
```

- network-overload-protection は、ネットワーク過負荷保護を示します
- mme-tx-msg-rate-control は MME メッセージ レート制御を示します
- egtp-pathfail は EGTP パス障害に関するメッセージ レート制御を示します

- ecm-idle は、ECM-Idle モードの MME UE セッションに関するレートを示します
- ecm-connected は、ECM-Connected モードの MME UE セッションに関するレートを示します
- <rate in sessions per second> は、セッションの 1 秒あたりのレートしきい値を 1 ~ 5000 の範囲で指定します

設定例

```
ddn throttle arp-watermark 1 rate-limit RATE time-factor 120 throttle-factor 50
increment-factor 10 poll-interval 30 throttle-time-sec 0 throttle-time-min 1
throttle-time-hour 0 stab-time-sec 0 stab-time-min 2 stab-time-hour 0
```

拡張輻輳制御

拡張輻輳制御機能を使用すると、MME は接続先の eNodeBs に信号を送信することで、MME プール内の他の MME にトラフィックをリダイレクトできます。これは S1 インターフェイス過負荷手順 (TS 36.300 および TS 36.413) によって実現します。

過負荷制御が設定されている場合、輻輳しきい値に達すると、MME に接続している一定割合の eNodeBs に S1AP インターフェイス過負荷開始メッセージを送信するよう、MME を設定することができます。MME で低減すべき負荷の量を反映するために、この割合は設定可能です。eNodeBs に送信される過負荷応答情報要素 (IE) の中で、MME は、特定タイプのセッションを拒否または許可するよう eNodeB に要求できます。たとえば：

- 非緊急セッションを拒否
- 新規セッションを拒否
- 緊急セッションを許可
- 高優先度セッションとモバイル終端サービスを許可
- 遅延許容アクセスを拒否

輻輳制御機能を使用すると、ポリシーとしきい値を設定し、高負荷状態に直面したときのシステムの対応方法を指定できます。輻輳制御機能はシステムを監視して、システムが高負荷のときにパフォーマンス低下につながる可能性のある状態を見つけます。通常、これらの状態は一時的なもので (高い CPU 使用率またはメモリ使用率など)、すぐに解決します。ただし、特定時間内に継続的に、または多数のこうした状況が発生すると、サブスクリバセッションを処理するシステム機能に影響を及ぼす可能性があります。輻輳制御は、このような状態を特定し、状況に対処するポリシーを起動するのに役立ちます。

輻輳状態のしきい値

- システム CPU 使用率
- システム サービス CPU 使用率 (Demux カード CPU 使用率)
- システム メモリ使用率
- ライセンス使用状況
- 1 サービスあたりの最大セッション数

しきい値と許容レベル

クリティカル、メジャー、マイナー輻輳レベルに関するしきい値と許容値を設定するときには、

しきい値レベルと許容値が決して重複することのないようにしてください。次に示す設定例では、しきい値レベルが重複しません。

- クリティカル輻輳は 95% でトリガーされ、90% でクリアされる
- メジャー輻輳は 90% でトリガーされ、85% でクリアされる
- マイナー輻輳は 85% でトリガーされ、80% でクリアされる

サービス コントロール CPU しきい値

このしきい値は、システムの demux CPU から計算されます。しきい値は、5 分間の CPU 平均使用率に基づいて計算されます。

demux CPU の 2 つの CPU コアの最大 CPU 使用率の値が考慮されます。たとえば、CPU コア 0 の 5 分間の CPU 使用率が 40% で、CPU コア 1 の 5 分間の CPU 使用率が 80% である場合は、CPU コア 1 がしきい値の計算対象となります。

システム CPU しきい値

このしきい値は、すべての CPU での 5 分間の CPU 平均使用率を使って計算されます (スタンバイ CPU と SMC CPU を除く)。

すべての CPU の 2 つの CPU コアのうち、最も高い使用率の値が考慮されます。

システム メモリしきい値

このしきい値は、すべての CPU の 5 分間の平均メモリ使用量を使って計算されます (スタンバイ CPU と SMC CPU を除く)。

輻輳アクション プロファイルの設定

輻輳アクション プロファイルは、該当するしきい値を超えた後に実行できる一連のアクションを定義します。

輻輳アクション プロファイルと輻輳制御ポリシーとの関連付け

それぞれの輻輳制御ポリシー (クリティカル、メジャー、マイナー) を、1 つの輻輳制御プロファイルに関連付ける必要があります。

過負荷制御の設定

MME で過負荷状態が検出された場合、指定した割合の eNodeBs にその状態を報告し、設定済みアクションを着信セッションに対して実行するよう、システムを設定することができます。

reject-new-sessions に加えて、次の過負荷アクションも使用できます。

- permit-emergency-sessions-and-mobile-terminated-services
- permit-high-priority-sessions-and-mobile-terminated-services
- reject-delay-tolerant-access
- reject-non-emergency-sessions

設定例の説明

これにより、輻輳制御機能が有効になります：

```
ddn throttle arp-watermark 1 rate-limit RATE time-factor 120 throttle-factor 50
increment-factor 10 poll-interval 30 throttle-time-sec 0 throttle-time-min 1
throttle-time-hour 0 stab-time-sec 0 stab-time-min 2 stab-time-hour 0
```

輻輳アクション プロファイル (クレティカル/メジャー/マイナー) の定義

```
ddn throttle arp-watermark 1 rate-limit RATE time-factor 120 throttle-factor 50
increment-factor 10 poll-interval 30 throttle-time-sec 0 throttle-time-min 1
throttle-time-hour 0 stab-time-sec 0 stab-time-min 2 stab-time-hour 0
```

輻輳ポリシーの適用

```
ddn throttle arp-watermark 1 rate-limit RATE time-factor 120 throttle-factor 50
increment-factor 10 poll-interval 30 throttle-time-sec 0 throttle-time-min 1
throttle-time-hour 0 stab-time-sec 0 stab-time-min 2 stab-time-hour 0
```

関連情報

- [Cisco ASR 5000 Mobility Management Entity Administration Guide](#)
- [テクニカル サポートとドキュメント – Cisco Systems](#)