
Calpont InfiniDB[®]

パフォーマンスチューニングガイド

Release 3.5.1
Document Version 3.5.1-1
December 2012



Copyright © 2012 Calpont Corporation. All rights reserved.

InfiniDB および Calpont 製品名は、Calpont の商標です。他社およびその製品への参照は、各社が所有する商標を使用しており、参照のみを目的としています。

この文書の情報は予告なしに変更される場合があります。この文書にいかなる誤りがある場合も、Calpont に責任はないものとします。

ソフトウェアライセンス

「GNU Free Documentation License」の条項のもと、この文書をコピー、配布、修正する権限が付与されます。Free Software Foundation によって発行されたバージョン 1.3 以降には、変更のない項、表紙、裏表紙が含まれていません。ライセンスのコピーについては「GNU Free Documentation License」の項に記載されています。

目次

	はじめに	i
	対象読者	i
	表記規則	ii
	マニュアルリスト	iii
	マニュアルの入手	iii
	マニュアルへのフィードバック	iii
	追加リソース	iii
第 1 章	パフォーマンスチューニングの概要	1
	InfiniDB のアーキテクチャの目標	1
	InfiniDB を使用した結合の高速化	2
	構文に関する注意点	2
	InfiniDB の主要コンポーネント	2
	分析用データベース InfiniDB のエディション	4
	InfiniDB のエディション間のチューニングの相違点	5
第 2 章	InfiniDB の分散処理モデル	7
	ユーザーモジュールからパフォーマンスモジュールへの リクエストの処理の粒度	7
	パフォーマンスモジュール内の処理の粒度（バッチ プリミティブ）	8
	バッチプリミティブステップ（BPS）	8
第 3 章	エクステントマップ	11
	スキャンによるエクステントマップ群	11
	他の列に対するパーティションブロックの除外	12
	列ストレージ、バッチプリミティブおよびエクステント マップ機能による I/O の除外	13

第 4 章	物理 I/O のチューニング	15
	最初のスキャン操作のチューニング	15
	追加の列の読取りのチューニング	17
第 5 章	同時実行および問合せ	19
第 6 章	InfiniDB の複数結合のチューニング	21
	簡単な 2 つの表の結合の詳細	21
	複数結合の詳細	22
	結合の最適化	23
	主要なチューニングパラメータ : PmMaxMemorySmallSide	24
	単一サーバーのインストールの一般的なチューニング ガイドライン	24
	複数 PM のインストールの一般的なチューニング ガイドライン	24
第 7 章	メモリー管理	27
	主要なチューニングパラメータ : NumBlocksPct および TotalUmMemory	27
第 8 章	スケーラビリティ	29
	スケーラビリティ : データのサイズとパフォーマンス	29
	スケーラビリティ : ユーザーモジュール	29
第 9 章	ツールおよびユーティリティのチューニング	31
	問合せのサマリー統計 : calgetstats	31
	問合せの詳細統計 : calgettrace	33
	キャッシュのフラッシュ : calflushcache	35
	パラメータの読取りまたは変更 : configxml.sh	36
	エクステンタマップの表示 : editem	36
第 10 章	列データストレージの相違点	37

第 11 章	データのロード速度およびリアルタイムに近いロード	39
	処理の優先順位付け	39
第 12 章	InfiniDB のパフォーマンスの目安	41
第 13 章	追加リソース、ダウンロードおよびサポート	45
付録 A	GNU Free Documentation License	47



はじめに

『Calpont InfiniDB パフォーマンスチューニングガイド』へようこそ。行ベースのデータベースでの経験が豊富な開発者または DBA の観点から見ると、InfiniDB 操作には従来のデータベース操作に関連するものと、直接関係のないものがあります。また、従来の行ベースの DBMS システムに共通の基礎となる操作が、InfiniDB 内には存在しない場合があります（たとえば、InfiniDB では全表スキャンは実行されません）。

本書は、I/O 要件を大幅に削減して大容量データの大幅な平行化および拡張を可能にする、列指向 RDBMS の InfiniDB のチューニングに役立ちます。

分析用データベース InfiniDB での操作は、マルチスレッド化および分散（オプション）されたシステムの大容量データのスキャン、結合、集計が優れたパフォーマンス特性で行われるように最適化されています。ここで使用される大容量データは、数億行から数千億行に渡る場合があります、また、数百 GB から数百 TB のデータになる場合があります。InfiniDB は、これらの数値をはるかに超える場合でも対応できるように設計されています。また、InfiniDB では、スケールが小さい場合でも優れたパフォーマンスが提供されます。

本書には、Calpont InfiniDB Enterprise にのみ有効な情報が含まれています。

対象読者

本書は、以下に示す様々な役割の読者を対象としています。

- ◆ データベース管理者
- ◆ アプリケーションおよび Web の開発者
- ◆ データの設計者
- ◆ システムおよび Web の管理者

表記規則

本書では、次の表記規則およびユーザーへの警告を使用しています。

表 1: 表記規則

項目	説明
太字	表示されたとおりに入力する文字。 例： getLogInfo と入力します この場合、 getLogInfo と入力します。
斜体	変数またはプレースホルダ。文字列を適切に置き換えて入力します。複数の単語で構成される変数はアンダースコア () で連結して表示されています。 例： <i>ID</i> を入力します。 ID 番号 34878 を入力します。 <i>IP_address</i> を入力します。 IP アドレス 110.68.52.01 を入力します。

表 2: ユーザーへの警告

項目	説明
	注意：役立つ情報であることを示します。
	警告：データの損失または破損の原因となるハードウェアやソフトウェアのエラーを発生させる可能性があることを示します。

マニュアルリスト

Calpont InfiniDB のマニュアルは、様々な読者を対象とした複数のガイドで構成されています。次の表を参照してください。

表 3: マニュアル

マニュアル	説明
『Calpont InfiniDB 管理者ガイド』	Calpont InfiniDB を管理するための詳細な手順について説明します。
『Calpont InfiniDB 最小推奨仕様ガイド』	Calpont InfiniDB の実装に必要なハードウェアおよびソフトウェアの最小の推奨仕様を示します。
『Calpont InfiniDB インストールレーションガイド』	分散構成に Calpont InfiniDB をインストールするために必要な手順の概要について説明します。
『Calpont InfiniDB SQL 構文ガイド』	Calpont InfiniDB に固有の構文について説明します。
『Calpont InfiniDB 概要』	分析用データベース Calpont InfiniDB の概要について説明します。
『Calpont InfiniDB マルチ UM 同期ガイド』	Calpont InfiniDB で 2 つ以上のユーザーモジュールの同期を保持するために使用するオプションの概要について説明します。

マニュアルの入手

英語版のマニュアルは、(<http://www.infinidb.org/> および <http://www.calpont.com>) で入手することができます。追加の支援が必要な場合は infinidb_doc@ashisuto.co.jp にご連絡ください。

マニュアルへのフィードバック

マニュアルの改善に向けて、フィードバック、コメントおよび提案をいただけますようお願いいたします。マニュアル名、バージョンおよびページ番号を添えてコメントを infinidb_doc@ashisuto.co.jp にご送付ください。

追加リソース

Calpont InfiniDB のインストールおよびチューニング、または Calpont InfiniDB を使用したデータの間合せに関して支援が必要な場合は infinidb_doc@ashisuto.co.jp までご連絡ください。

パフォーマンスチューニングの概要

InfiniDB のアーキテクチャの目標

パフォーマンスに関連する設計を決定づける多くのコアアーキテクチャの目標があります。InfiniDB の目標は次のとおりです。

1. メモリアクセスおよびストレージアクセスの両方で問合せの I/O コストを大幅に削減する。
 - 大規模データに対するほとんどの問合せで必要な I/O を根本的に削減する。
 - データのブロック（ページ）へのランダムなアクセスを排除する。
 - 大規模なデータセットの部分キャッシュまたは完全キャッシュを可能にするグローバルなデータバッファキャッシュを実装する。
2. 最小限の同期で、データベース操作を並行して実行できるようにする。
 - データのブロックの読取りに関連する同期を排除する。
 - 重要なスレッドプールに対して送信および受信を行うキューなど、多くのメカニズムを使用してスレッドの同期の問題を最小限にする。
3. 次のようなパラレルのデータベース操作の作業単位を定義する。
 - 様々なストレージシステムで正常に機能する。
 - パフォーマンスモジュールが常に 100% 近くの CPU 使用率で動作できるようにする。
 - スレッドを問合せ専用にしなない（つまり、ロード時でも小規模の問合せが短時間で実行されるようにする）。
4. データの送信コストを最小限にする（つまり、ローカルスレッド間または分散されたスレッド間でのデータの送信を最小限にする）。
 - 可能な場合、データを送信するのではなく、データの操作を送信する。
 - 可能な場合、最大の表に対する送信コストが発生しないように操作を設定する。

5. チューニング要件を最小限にし、詳細データの非定型分析のパフォーマンスが向上する。

現時点では、1行の検索にかかる経過時間の最小化は主要な目標に含まれないことに注意してください。10億行から1行を検索するパフォーマンスは一瞬にして実現できますが、従来の索引と表の組合せほど効率的ではない場合があります。

InfiniDB を使用した結合の高速化

InfiniDB では、従来の MySQL の結合を使用せずに、InfiniDB エンジン内で結合が実現されます。これらの結合の特性は次のとおりです。

1. ハッシュ結合操作によって数百万または数十億の行に対して最適化されます。
2. スケーラブルなスレッドプールでマルチスレッド化および分散されます。
3. 中間結果をマテリアライズせずに任意の数のディメンション表に対して大規模なファクト表を1回のパスでストリームします。
4. データ送信コストを大幅に削減する可能性がある集計操作を行って、送信されません。

このマニュアル内のパフォーマンスチューニングのガイドラインは、分析用データベース InfiniDB 内でのみ実行される結合およびその他の操作に適用されます。

構文に関する注意点

InfiniDB では、使用可能なすべての構文に対して完全なパフォーマンス機能を備えているわけではありません。サポートされている構文については、『Calpont InfiniDB SQL 構文ガイド』マニュアルを参照してください。新しい機能が追加されると、構文に対するその拡張機能では、他の操作の場合と同様に、マルチスレッド化されたスケラブルなパフォーマンスが積極的に活用されます。

現在、追加の構文は代替構成（つまり、InfiniDB を標準ストレージエンジンとして構成すること）によって利用できます。ただし、この場合は、マルチスレッド結合または分散結合および集計操作を行うことができなくなります。

InfiniDB の主要コンポーネント

InfiniDB では、3つの主要コンポーネントで構成されるモジュールアーキテクチャが提供されます。3つすべてのコンポーネントが連携して InfiniDB インスタンスを構成しています。これらのコンポーネントには次のものがあります。

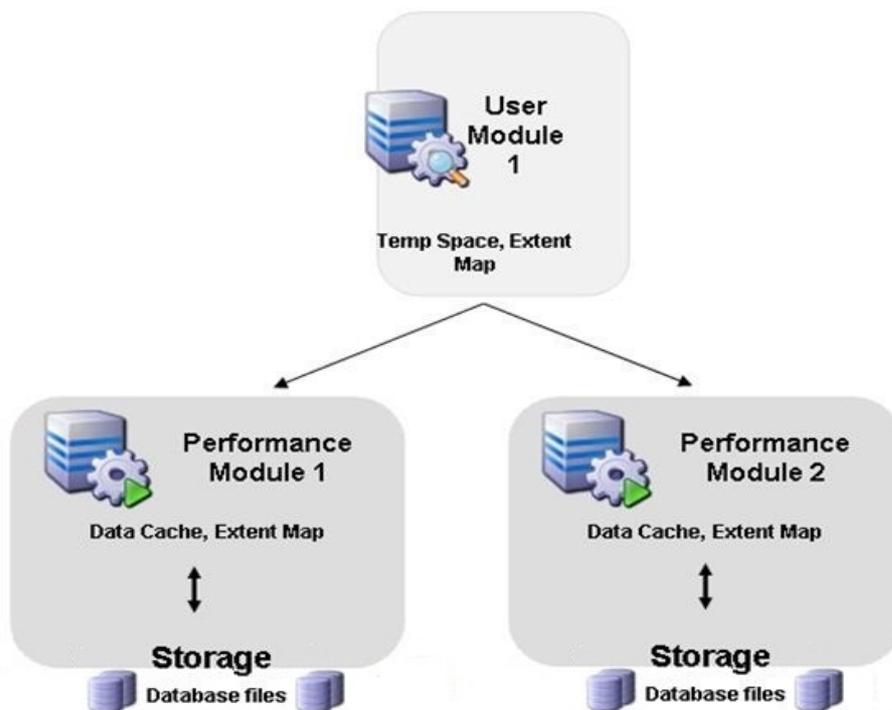
- ◆ **ユーザーモジュール (UM)** : ユーザーモジュールは、SQL リクエストを分割して1つ以上のパフォーマンスモジュールに分散します。パフォーマンスモジュールに

よって、リクエストされたデータがメモリーキャッシュまたはディスクのいずれかから実際に取得されます。1つのユーザーモジュールで1つの問合せの状態が保持されます。

- ◆ **パフォーマンスモジュール (PM)** : パフォーマンスモジュールは、マルチスレッド方式で、ユーザーモジュールから受信した細かい作業単位を実行します。また、InfiniDB Enterprise では、多数のパフォーマンスモジュール間で作業を分散できます。
- ◆ **ストレージ** : InfiniDB では、ローカルストレージまたは共有ストレージ (SAN など) のいずれかを使用してデータを格納できます。ユーザーは、1つのみのサーバーを InfiniDB サーバーとして動作させ、すべてを構成してそのサーバー上で実行するか、または複数のサーバーにスケールアウトできます。

また、InfiniDB では、エクステントマップと呼ばれる共有オブジェクト内の各列に関するメタデータが保持されます。エクステントマップは、ユーザーモジュールで発行する操作とその対象のデータを判断するときに参照されます。さらに、ディスクからのブロックの読取りに必要な場合はパフォーマンスモジュールによって参照されます。各列は1つ以上のファイルで構成され、各ファイルに複数のエクステント (通常はデータの連続した割当て) を含めることができます。エクステントも、エクステントマップ内で追跡されます。

UM では問合せ計画が認識され、PM ではデータブロックおよび操作が認識されます。エクステントマップによってこの抽象化が可能になります。



2つのパフォーマンスモジュールがインストールされた環境を簡単に示した図

分析用データベース InfiniDB のエディション

InfiniDB は、大規模なデータのロードおよび問合せを行うために構築された列指向のデータベースで、オープンソース版および製品版の両方で使用できます。InfiniDB のすべてのエディションに、パフォーマンスで重要となる次の機能が含まれています（機能についての詳細は、他のマニュアルを参照してください）。

- ◆ **列指向のアーキテクチャ**：InfiniDB では、行ごとにではなく列ごとにデータが格納されます。これによって、スキャン操作でスキャンに含まれない列を無視したり、問合せで参照されない列を無視する列操作を選択することができます。
- ◆ **マルチスレッド設計**：InfiniDB では、問合せをサポートするために操作が非常に細かいレベルで分散され、同期が最小限または行われない状態でそれらの操作を実行できるようにほとんどの問合せが構築されます。操作はすべてのパフォーマンスモジュールに自動的にマップされ、分散されます。結果は削減され、コール側のユーザーモジュールに返されます。

- ◆ **垂直および水平の自動パーティション化**：列ストレージに基づいた I/O の大幅な削減に加えて、InfiniDB では、ほとんどの列に対して小さいエクステンツ（パーティション）が作成され、状況によってはエクステンツの除外を発生させることができるように一部のメタデータが格納されます。
- ◆ **柔軟な同時実行**：InfiniDB では、データへの同時アクセスが可能です。これによって、大規模な問合せで、使用可能なすべてのパフォーマンスモジュールのリソースを短時間使用できるだけでなく、大規模な問合せが完了するまで待機せずに小規模な問合せを実行できます。

InfiniDB Enterprise では、製品版ライセンスが使用され、次のパフォーマンス関連の機能を含むスケールアウトオプションが追加されます（機能についての詳細は、他のマニュアルを参照してください）。

- ◆ **大規模パラレル処理（MPP）対応**：InfiniDB は複数の汎用ハードウェアマシンを使用して、総合的なパフォーマンスをほぼリニアに向上させることができます。複数のパフォーマンスモジュールを使用すると、使用可能なすべてのパフォーマンスモジュール上のすべてのスレッド間で細かい作業を分散できます。
- ◆ **シェアードナッシング分散データキャッシュ**：複数ノードによる InfiniDB 構成では、データは様々なノードとそのデータキャッシュに分散されます。ノード間でデータを共有することはありませんが、問合せのためにデータを読み取る時は、InfiniDB MPP アーキテクチャ内のすべてのノードがアクセスされます。つまり、InfiniDB は、すべての参加ノードが並行して分散形式でアクセスする、1つの大きな論理データキャッシュを作成すると考えられます。これによって、十分な数のノードがあれば、InfiniDB は文字通り、大規模なデータベースをキャッシュできます。

InfiniDB のエディション間のチューニングの相違点

このチューニングガイドでは、オープンソース版およびエンタープライズ版の両方について説明しています。エンタープライズ版に限定されたいくつかの追加のチューニング概念について特に詳しく説明しています。また、いくつかのデータベースを最初に単一サーバーとしてインストールした後、パフォーマンス要件またはスケール要件を満たすために、または分散されたシステムの管理を簡単にする追加の監視機能および管理機能を活用するためにエンタープライズ版にアップグレードすることをお勧めします。

InfiniDB の分散処理モデル

InfiniDB のジョブ発行処理は、Map/Reduce 処理に例えることができます。UM は、スレッドプール（分散される場合がある）に操作を発行します。このスレッドプールは、データに対して個別にデータベース操作（フィルタ、集計、結合など）を実行し、次の操作を行うためにデータセット（削減されている場合がある）を UM に戻します。

ただし、InfiniDB と Map/Reduce にはいくつかの相違点があります。ジョブ（プリミティブ操作と呼ばれる）は SQL 構文にマップ済で、利用可能な外部 API はないため、これはツールキットやソフトウェアコールではありません。また、ジョブのスケジューリング処理にも相違点があります。

ユーザーモジュールからパフォーマンスモジュールへのリクエストの処理の粒度

データベース操作のリクエストはユーザーモジュールから 1 つ以上のパフォーマンスモジュールに発行され、800 万行を含むデータブロックの範囲（エクステントと呼ばれる）に対してスキャンが実行されます。データがまだキャッシュされていない場合は、リクエストにより、この操作に加えてストレージからの読取りが発生する場合があります。システムのマルチブロック読取り機能を最大限に活用するために、1 つのエクステント内のすべてのブロックがターゲット PM にマップされ、発行されます。また、同じエクステントのデータを必要とするその他の問合せも同じ PM に送信されます。エンタープライズ版で PM の数が増加した場合は、新たに追加された PM を新しいリクエストに含めることができるように再マッピングが迅速に実行されます。非分散システムの場合でも、この粒度で操作を実行すると、隣接するブロックがまとめて読み取られて短い時間で処理されるなど、ブロックアクセスの空間的および時間的局所性が向上します。

大きい表内の部分的に移入されているエクステントと、小さい表用の部分的かつ唯一のエクステントの場合、操作のリクエストは各エクステント内の最高水位標（HWM）までに制限されることに注意してください。

エクステントのテスト済のデフォルト構成は 800 万行です。ディスク上のこの割当てサイズは、8-64MB（列内の 1 エクステントのディスク上のデータサイズ）の範囲で変わります。値はシステムのインストール時に設定されます。この値は、特定の InfiniDB インスタンス全体に対するものであり、スタートアップ後は変更できないことに注意してください。新しいインスタンスの動作は、別のマニュアルに記載されて

いる Calpont.xml パラメータによって設定できますが、追加の検証および確認が必要になります。

パフォーマンスモジュール内の処理の粒度 (バッチプリミティブ)

パフォーマンスモジュール内のデータベース操作は、よりきめ細かい方法で分割され、個々のスレッドはエクステント内の個々のデータブロックで個別に動作します。最小の操作はプリミティブとも呼ばれます。プリミティブ処理を行う各 PM 上で実際に実行されるソフトウェアのプロセスは PrimProc です。バッチプリミティブでは、区切られたブロックセット内に格納されているすべての行に対して 1 つ以上のデータベースプリミティブが実行されます。

InfiniDB のブロックは 8192 バイトのデータで、固定長ストレージにおいては 1024-8192 行をサポートし、varchar 列においては長さによって行数が変わります。varchar(7) およびそれより小さいものは固定長のフィールドにマップされ、固定長のフィールドとして格納されることに注意してください。

様々な状況において、ある特定のブロックがバッチプリミティブに含まれるかどうかを考える場合には、少し概念的な説明が役立ちます。実際、バッチプリミティブはデータサイズに応じて、1-8 個のブロックに格納されている 1 列の小さい固定された範囲の行 (8000 行) に対して操作を実行します。また、バッチプリミティブ内の複数の個々のプリミティブは、複数の列でこの小さい固定されたサイズ範囲の行に対して操作を実行することがあります。そのため、最小のプリミティブは単一の 8K ブロックのデータを読み込むことで終了する場合がありますし、その他の場合は多数の列にわたる問合せに対応するために 100 を超える数の個々のブロックからの読取りを行う場合があります。

バッチプリミティブステップ (BPS)

バッチプリミティブステップ (BPS) は、操作が UM から 1 つ以上の PM に発行される問合せ実行計画ステップで、UM で実行される可能性がある補足ステップも含まれます。利用可能なすべての PM 内の個々のスレッドは、割り当てられた範囲の行に対して、リクエストされたバッチプリミティブを実行します。実際、サーバーが 1 つか多数かにかかわらず、バッチプリミティブの処理に利用できるグローバルスレッドプールがあります。

問合せは、1 つ以上のバッチプリミティブステップによって処理されます。バッチプリミティブステップは、問合せ要件に応じて、次の一部またはすべての処理を実行できます。

- ◆ **単一系列のスキャン**：単一系列の述語 (=、<>、in (list)、between、is null など) に基づいて、特定の列の 1 つ以上のエクステントをスキャンします。
詳細は、15 ページの「最初のスキャン操作のチューニング」を参照してください。

- ◆ **単一系列の追加フィルタ**：以前のスキャンで検出された行に対する追加の列を投影し、必要に応じて単一系列の追加述語を適用します。ブロックのアクセスは行識別子に基づいており、ブロックに直接アクセスされます。

詳細は、17 ページの「追加の列の読取りのチューニング」を参照してください。

- ◆ **表レベルのフィルタ**：表レベルの任意のフィルタ（`column1 < column2` や、より高度な関数および式など）に必要な追加の列を投影します。この場合も、ブロックのアクセスは行識別子に基づいており、ブロックに直接アクセスされます。

- ◆ **結合での結合列の投影**：任意の結合操作に必要な追加の結合列を投影します。この場合も、ブロックのアクセスは行識別子に基づいており、ブロックに直接アクセスされます。

詳細は、21 ページの「InfiniDB の複数結合のチューニング」を参照してください。

- ◆ **複数結合の実行**：投影された結合列に対して 1 つ以上のハッシュ結合操作を適用し、その値を使用して、以前に構築されたハッシュマップを調べます。内部結合または外部結合の要件を満たすのに必要なタプルを構築します。エンタープライズ版では、以前に構築されたハッシュマップのサイズに応じて、PM 処理を実行しているサーバーまたは UM 処理を実行しているサーバーで実際の結合動作を実行できることに注意してください。いずれの場合も、バッチプリミティブステップの機能は同じです。

詳細は、21 ページの「InfiniDB の複数結合のチューニング」および 27 ページの「メモリー管理」を参照してください。

- ◆ **クロス表レベルのフィルタ**：クロス表レベルの任意のフィルタ（`table1.column1 < table2.column2` や、より高度な関数および式など）に必要な追加の列をプリミティブステップの行の範囲から投影します。この場合も、ブロックのアクセスは行識別子に基づいており、ブロックに直接アクセスされます。前提条件の結合操作が UM で実行された場合は、この操作も UM で実行されます。それ以外の場合は、PM で実行されます。

- ◆ **集計 / 重複行の削除操作 (1 回目)**：特定のバッチプリミティブに割り当てられた結合済の行のセットに対してローカルのグループ化操作、重複行の削除操作または集計操作を適用します。この処理の 1 回目は、パフォーマンスモジュールによって処理されます。

- ◆ **集計 / 重複行の削除操作 (2 回目)** : 特定のバッチプリミティブに割り当てられた結合済の行のセットに対して最後のグループ化操作、重複行の削除操作または集計操作を適用します。この処理は、ユーザーモジュールによって処理されます。
集計の詳細は、27 ページの「メモリー管理」を参照してください。

バッチプリミティブは、I/O を最小限に抑えるために、フィルタをできるかぎり早く適用し、追加の列の投影をできるかぎり遅くすることによって、列データファイルに対して表指向の SQL コマンドを効率的に実行します。さらに、バッチプリミティブは、グループ化操作、集計操作、重複行の削除操作を実行して、ユーザーモジュールに戻されるバイト数を削減します。

エクステントマップ

エクステントは、物理的なセグメントファイル内に存在する、指定した列の領域の論理ブロックで、サイズは 8-64MB の間です。各エクステントでは、デフォルトで 800 万行がサポートされており、小さいデータ型であれば、ディスク上で消費される領域は少なくなります。

エクステントマップは、ストレージに保持されるすべてのデータのカタログとみなすことができます。各エクステントに 1 つのエントリが含まれます。各カタログエントリには、ブロックの範囲ごとの論理識別子（行識別子の一部）、部分的にデータが移入されている各エクステントの最高水位標、ほとんどのデータ型に対する各エクステントの最小値および最大値を保持する場所が含まれています。現在、8 バイトを超える文字データ型および 7 バイトを超える varchar データ型では最小値および最大値は移入されません。日付型、10 進型、整数型、小さい文字列などの他のすべてのデータ型では最小値および最大値が移入されます。

スキャンによるエクステントマップ群

エクステントマップは DML またはバルクロード機能によって自動的に移入されるため、特別な統計収集は必要ありません。これらの最小値および最大値はディスクに保持されるため、システムの再起動後でもその情報を使用できます。

列の各エクステントに格納された最小値および最大値は、特定の状況で I/O を除外するために使用され、セッションで最近実行された文に関する情報を表示する `select calgetstats()` ; 関数を使用して表示できます。次の例では、BlocksTouched および PartitionBlocksEliminated に対してレポートされる値を示します。

BlocksTouched : 問合せをサポートするためにアクセスされた 8KB のデータブロックの数。

PartitionBlocksEliminated : 最小値と最大値の比較にのみ基づいて除外されるブロックの数。

次の問合せは、60 億行が含まれているスケール係数 1000 のスタースキーマベンチマーク (SSB) データセットに対して実行されます。lo_orderdate フィールドに基づいて一度に 1 か月分のデータがロードされており、ソートは実行されていません。この問合せでは、(大幅に値下げされた項目を除外するために) 7700 万のレコードに関して価格とコストを比較して、特定の日付範囲の分析を行います。その後、約 6000 万の個別の計算値の平均が計算されて、純収入が特定されます。

```

select lo_discount, avg(lo_extendedprice - lo_supplycost),
count(*)
  from lineorder
  where lo_orderdate between 19940101 and 19940131
        and lo_supplycost < lo_extendedprice * .5
group by 1 order by 2;

```

```

+-----+-----+-----+
| lo_discount | avg(lo_extendedprice - lo_supplycost) | count(*) |
+-----+-----+-----+
|          3.00 |                3806802.786682 | 5510648 |
|          5.00 |                3807880.542760 | 5511323 |
|          6.00 |                3808025.905069 | 5511063 |
|          9.00 |                3808060.552264 | 5503826 |
|          7.00 |                3808450.935915 | 5506420 |
|         10.00 |                3808529.811736 | 5508111 |
|          1.00 |                3808842.678177 | 5509157 |
|          0.00 |                3809090.349685 | 5510113 |
|          2.00 |                3809309.979333 | 5507300 |
|          4.00 |                3809375.700146 | 5506796 |
|          8.00 |                3810136.417883 | 5509236 |
+-----+-----+-----+
11 rows in set (3.20 sec)

```

```

+-----+
| calgettats()
|
+-----+
| Query Stats: MaxMemPct-0; NumTempFiles-0; TempFileSpace-0MB; ApproxPhyI/O-0;
CacheI/O-456180; BlocksTouched-456180; PartitionBlocksEliminated-
2890042; MsgBytesIn-4MB; MsgBytesOut-0MB; Mode-Distributed| 1256230262 897059 |
+-----+
1 row in set (0.00 sec)

```

他の列に対するパーティションブロックの除外

最小値または最大値に基づいたパーティションブロックの除外は他の列に対しても発生する場合があります。通常は注文日に関連する別の日付フィールド（`lo_commitdate`）を使用するように問合せを変更すると、ブロックアクセスと経過時間の両方で同様の削減が行われます。昇順のキー値などの1つ以上のデータパターンを含むほとんどのデータセットで、同様のメリットが見込まれます。

```

mysql> select lo_discount, avg(lo_extendedprice - lo_supplycost), count(*)
->   from lineorder
->   where lo_commitdate between 19940101 and 19940131
->   and lo_supplycost < lo_extendedprice * .5

```

```
-> group by 1 order by 2 ;
```

```
+-----+-----+-----+
| lo_discount | avg(lo_extendedprice - lo_supplycost) | count(*) |
+-----+-----+-----+
|          1.00 |                3807548.123096 | 5509631 |
|          2.00 |                3807764.108009 | 5506100 |
|          3.00 |                3808203.553947 | 5508280 |
|         10.00 |                3808530.528851 | 5508275 |
|          0.00 |                3808644.366889 | 5508041 |
|          8.00 |                3808727.418191 | 5507120 |
|          6.00 |                3809173.101862 | 5507031 |
|          4.00 |                3809189.422454 | 5508216 |
|          9.00 |                3809475.131073 | 5507767 |
|          5.00 |                3809746.089397 | 5511763 |
|          7.00 |                3809870.163028 | 5510417 |
+-----+-----+-----+
```

11 rows in set (4.08 sec)

```
mysql>
mysql> select calgetstats();
+-----+
| calgettats() |
+-----+
+-----+
| Query Stats: MaxMemPct-0; NumTempFiles-0; TempFileSpace-0MB; ApproxPhyI/O-0;
CacheI/O-1561569; BlocksTouched-1561569; PartitionBlocksEliminated-
2779450; MsgBytesIn-14MB; MsgBytesOut-0MB; Mode-Distributed| 1256231405 276729
|
+-----+
1 row in set (0.00 sec)
```

列ストレージ、バッチプリミティブおよびエクステントマップ機能による I/O の除外

InfiniDB の列ストレージでは、問合せに含まれていない列に対する I/O を自動的に削減できます。また、バッチプリミティブステップ内のデータをフィルタ処理する操作を積極的に使用することで、フィルタに従って参照される列の I/O を大幅に除外できます。さらに、InfiniDB では、バッチプリミティブの発行前に問合せが分析されるため、各列に格納されている最小値および最大値を利用することで、最初のスキャンの I/O も削減される可能性があります。

それぞれの I/O の除外方法を図で示すことができます。5つの列と、列ごとに13個のエクステント（12個は完全に移入され、1つは半分だけ移入されている）に分散された1億の行が含まれている単純な表があるとします。列 a および b のフィルタに基づいて表から列 a、b および c を選択する問合せでは、問合せを処理するためにアクセスされるエクステントは65個のうち5つのみとなる可能性があります。

◆ 列ストレージの最適化

- 列 d および e のブロックは、問合せで参照されなかったため無視されます。次の図では、除外されたエクステントは黄色で表示されています。

Extent

◆ エクステントマップの最適化

- 列 a のフィルタに対するエクステントマップの最小値および最大値は、エクステント 1-9 を除外します。次の図では、除外されたエクステントは緑色で表示されています。

Extent

- 列 b のフィルタはエクステント 9-11 を除外します。これによって、エクステント 12 および 13 のみがスキャンされます。次の図では、除外されたエクステントは青色で表示されています。

Extent

◆ バッチプリミティブの最適化

- 列 a および b に対してフィルタを実際に行うことによってエクステント 13 内の行が除外され、データセットがエクステント 12 内の数行に絞られた場合、参照されるのは列 c のエクステント 12 とブロックの小さいサブセットのみが参照されます。次の図では、列 c の除外されたエクステントはオレンジ色で表示されています。

Extent

Extent #	column a	column b	column c	column d	column e
1	Extent	Extent	Extent	Extent	Extent
2	Extent	Extent	Extent	Extent	Extent
3	Extent	Extent	Extent	Extent	Extent
4	Extent	Extent	Extent	Extent	Extent
5	Extent	Extent	Extent	Extent	Extent
6	Extent	Extent	Extent	Extent	Extent
7	Extent	Extent	Extent	Extent	Extent
8	Extent	Extent	Extent	Extent	Extent
9	Extent	Extent	Extent	Extent	Extent
10	Extent	Extent	Extent	Extent	Extent
11	Extent	Extent	Extent	Extent	Extent
12	Extent	Extent	Extent	Extent	Extent
13	Extent	Extent	Extent	Extent	Extent

1 億行、5 列のエクステントの図。

物理 I/O のチューニング

InfiniDB で物理 I/O をチューニングする最良の方法について説明する前に、まずこの項で使用するいくつかの用語を定義します。

- ◆ **エクステント**：通常、エクステントはストレージシステム内で連続して割り当てられた領域です。表の増大とともに、必要に応じて各列でエクステントが追加されます。エクステントは、ほとんどの場合、実際には 8-64MB の範囲の可変サイズで、800 万行を格納します。
- ◆ **DBRoot**：InfiniDB インスタンスの作成時に InfiniDB インスタンスで使用可能になるマウントポイント。
- ◆ **マルチブロック読取り**：ストレージからの多数の個々のブロック読取りを 1 つの読取り操作にまとめます。ほとんどのストレージサブシステムでは、主としてリクエストされた読取り操作のサイズに応じて、異なる持続速度（帯域幅）でデータブロックが提供されます。個々のランダムなブロックがリクエストされると、各ブロックをサポートするために機械的なディスクヘッド移動が発生する可能性があるため、ディスクパフォーマンスが大幅に低下することがあります。

データベース内の特定の表に行が追加されると、必要に応じてエクステントが追加され、その表のすべての列のサイズが増大します。特定の表の各列の最初のエクステントセットが 1 つの DBRoot に書き込まれます。同じ表の 2 番目のエクステントセットはラウンドロビン方式で次の DBRoot に書き込まれます。その他の表も異なる DBRoot で始まります。このデータの分散は、列ファイルの同時実行読取りの様々な条件を満たすために使用可能なストレージリソースを最大限に活用するために設計されています。

前述の動作は、定義されたエクステントサイズに基づいてデータがロードされると自動的に発生します。実際に、表作成操作に関連する追加のストレージパラメータは存在しません。

最初のスキャン操作のチューニング

I/O をチューニングするには、バッチプリミティブステップ（BPS）で実行される操作の順序を理解する必要があります。0 個以上のフィルタを使用して、特定の BPS に対して 1 つ以上の列が参照されます。すべての場合で、まず 1 つの列が読み取られ、この最初の操作をサポートするためにエクステント全体が読み取られます。前のフィルタの選択に応じて、追加の列ブロックが読み取られる場合があります。ただし、追

加の列では、個々のブロックアクセスを可能にする行識別子を使用できます（指定されている場合）。

「最初のスキャン」操作には、エクステント内のすべてのブロックが必要です。「追加の列の読取り」は、すべてのブロックのスキャン（マルチブロック読取りを最大活用）または個々のブロックの読取り（一部の状況でのブロックの合計数の最小化）のいずれかを行うことで実現できます。

通常、最初のスキャン操作のチューニングは、2つの相反する目標および主要作業単位（エクステント）に基づいて行われます。目標は次のとおりです。

1. 作業単位（エクステント）に対して並行して処理を行うためにできるかぎり多くのスレッドを使用可能にする。読取りは本質的に非同期であるため、スレッドはI/Oの待機中に停止されません。
2. ストレージ構成を最大限に活用するためにマルチブロック読取りを利用する。

InfiniDBは多くの異なるサーバー構成またはストレージ構成で実行されることが予測されるため、これらのパラメータを調整すると所定の環境のパフォーマンスに影響を及ぼす可能性があります。デフォルト設定によって、様々な構成で高パフォーマンスが得られる必要があります。

この動作を制御する Calpont.xml 構成ファイル内のパラメータを次に示します。

```
ColScanReadAheadBlocks
```

デフォルト設定は、512（データのブロック数）です。

使用例：8バイトのデータ型に、スキャン操作をサポートするために読み取られる8192個のデータのブロックが含まれています。一度に128個以上のブロックを読み取るときにスループットが最大になるストレージサブシステムでは、ColScanReadAheadBlocksパラメータに128、256または512以上の値を指定すると、ストレージサブシステムの最良のパフォーマンスを実現できます。パラレル読取りの実行に8つのコアを使用可能なサーバーでは、1024、512または256以下の値を指定すると、（最大で）8つすべてのコアで並行して読取りを実行できます。

多くのサーバー構成またはディスク構成では、類似する結果をもたらす有効な設定が複数存在する可能性があります（これらの設定では、次の両方の要件が満たされるため）。

1. 並行して読み取るのに十分なスレッドを使用できる。
2. 帯域幅を最大化するためにマルチブロック読取りを十分に利用する。

追加の列の読取りのチューニング

追加の列の読取りに対する物理 I/O のパフォーマンスは、前の列に適用されたフィルタの特定のカーディナリティまたはフィルタチェーンによって異なります。エクステントから 1 つの行のみが必要な場合は、どのストレージシステムでも、そのブロックのみを読み取ると時間が短縮されます。いくつかのブロック、ほとんどのブロックまたはすべてのブロックが必要な場合は、一部の個々のブロックが問合せのサポートで参照されない場合でも、マルチブロック読取りを利用すると時間が短縮される可能性があります。

追加の列の読取り操作で使用可能な読取り方法は 2 つあります。

1. 既知のアドレスに基づいて個々のブロックを読み取る。
2. `ColScanReadAheadBlocks` に設定されている数のブロックを一度に読み取るマルチブロック読取りを使用する。

読取りが必要な場合、これら 2 つの操作の選択は PM 内で行われます。実際には、同じ列に対する以前の読取りから統計が収集され、パラメータと比較されます。

この動作を制御する `Calpont.xml` ファイル内のパラメータを次に示します。

`PrefetchThreshold`

デフォルト設定は、5（問合せで使用されたブロックの割合）です。

デフォルト設定の 5 は、以前読み取られたブロックの 5% を超えるブロックが問合せを満たすために実際に必要だった場合に、マルチブロック読取りが継続して使用されることを示します。問合せを満たすために必要なブロックの割合が 5% を下回った場合は、後続の読取りで個々のブロックの検索が使用され始めます。

このパラメータは、行の数ではなく、フィルタで必要なブロックの数に関連していることに注意してください。行の 5% を返す述語には、90% を超えるブロックが必要な場合があります。ただし、データ値が高度にクラスタ化されている場合、ブロックの数は 5-10% まで下がる場合があります。

実際の動作は、多くの項目の複雑な相互関係に関連しています。

1. 問合せで使用される特定のフィルタまたはフィルタチェーン。
2. フィルタで使用される列の値の頻度。
3. ブロック内でのこれらの値の実際の分布。
4. その時点でのデータバッファキャッシュ内のブロック。

5. 個々のブロックを読み取った場合と総ブロックを読み取った場合のストレージシステムの相対コスト。
6. ストレージシステムに対する同時実行問合せの影響。

通常、このパラメータはほとんどの問合せに影響を及ぼしません。これは、多くの InfiniDB の機能および動作に基づいています。

1. 複数の最適化に基づいた I/O の大幅な削減によって、すべての状況でストレージシステムのボトルネックの発生を防止できます。
2. 多くの問合せで I/O 要件が軽減されるため、より多くの問合せをキャッシュによって処理できます。また、エンタープライズ版では、グローバルなスケラブルデータキャッシュを使用できます。
3. 通常、ほとんどの分析的な問合せは、エクステント内の多くのブロックにアクセスし、境界条件に該当しません。

同時実行および問合せ

InfiniDB では、同時実行の管理は、特定の問合せに割り当てるスレッドの数を調整したり、スレッドの優先順位を割り当てるのではなく、バッチプリミティブ操作のリクエストが UM から PM に発行される割合を管理することによって行われます。

処理の流れは次のパラメータによって制御されます。

`MaxOutstandingRequests`

デフォルト値は 5 (エクステント) です。

前述のとおり、UM から発行されるバッチプリミティブリクエストは、1 エクステント (各列に 1 つのエクステント) 内の行の範囲に対して特定の操作を実行するリクエストです。この処理では、次の一般的なフィードバックループが行われます。

- ◆ UM は、`MaxOutstandingRequests` に設定されている数まで、バッチプリミティブを発行します。
- ◆ PM は、データブロックを処理し、個々の応答を戻します。
- ◆ UM は、すべての未処理バッチプリミティブからのブロック応答を個別に受信し、次のエクステントに対して次のバッチプリミティブを発行するタイミングを決定します。

たとえば、パラメータで 5 (エクステント) を設定した場合、5 つのエクステントを処理するバッチプリミティブリクエストが発行されます。ブロック応答が戻されると、UM は、未処理リクエストの総数がパラメータの設定に達するように、新しいエクステントに対して 1 つの追加バッチプリミティブを発行します。戻されたブロックのサイズの合計が 1 エクステントのサイズに達すると、次のエクステント用に別のバッチプリミティブが発行されます。

このパラメータの設定では、前述のパラメータ設定を使用してどの時点でも 5 つ相当のエクステントがアクティブに処理されようにすることが目標として定義されます。

事実上、この処理によって、利用可能なすべてのリソースを大規模な問合せで使用でき (他で消費されていない場合)、小規模な問合せを最小限の遅延で実行できます。各列の 1 つのエクステントに対して個々のバッチプリミティブを実行するのに必要な

時間は非常に短く、ディスクから多数の列を読み取る場合にかかる時間は最大で1秒ほどです。

負荷が大きい場合に小規模な問合せに対する遅延を最小限に抑えるには、小さい整数値を設定することをお勧めします。

大規模な問合せを優先させるには、MaxOutstandingRequests の値を少し大きくします。完全または完全に近い CPU 使用率で PM が動作できる状態にシステムが安定するまで、少しだけ時間がかかる場合があります (1秒ほど)。多くの場合、デフォルト値によってすべての PM が常に動作している状態に保たれるため、キュー内の追加リクエストによって経過時間が変わることはありません (開始時間は短縮されません)。

通常、InfiniDB Enterprise では、単独で実行される大規模な問合せが迅速に開始されて利用可能なすべてのリソースが使用されるように、パフォーマンスモジュールの数を基に MaxOutstandingRequests の設定を調整する必要があります。同時実行のワークロードシナリオでは、複数問合せのサポートによってキューが一杯になるため、このパラメータの設定にかかわらず、通常は完全なシステムの使用率が実現します。

システムのテストおよびベンチマークの結果では、単独で実行される大規模な問合せを迅速に開始できる各スケールアウト構成に対する適切な開始値は次のとおりです。

<u>PMs</u>	<u>MaxOutstandingRequests</u>
1-2	3
3-4	5
5-6	7
7-8	9
9 or more	PM count + 1

InfiniDB の複数結合のチューニング

InfiniDB の結合処理モデルでは、ネステッドループ操作は行われず、かわりにハッシュ結合操作が実行されます。比較的少ない数の行の結合（索引によってサポートされている場合）では、ハッシュ構造を構築する必要がないため、ネステッドループ操作の方が高速になる可能性が高くなります。一方、ハッシュ結合は、数千以上の行を結合する場合に一般的にパフォーマンスが高く、数百万、数十億という行を結合する場合はネステッドループ操作よりもかなり高速に実行されます。

InfiniDB エンジンには、ハッシュ結合操作を構築する多数の方法があります。これらの方法では、結合操作に関連する最大の表のハッシュ構造を作成する必要がないように、操作が優先的に順序付けられます。かわりに、最大の表（ファクト表）をマテリアライズする必要がないように、このファクト表には必要なすべてのフィルタ、結合および集計操作が行われます。ファクト表のマテリアライズを回避するような問合せの構築は、InfiniDB の最適化によって自動的に行われます。

InfiniDB が実行するハッシュ結合は、UM が管理するバッチプリミティブステップの範囲内で実行されます（したがって、PM が実行するバッチプリミティブ操作によって実行されます）。バッチプリミティブステップは問合せの実行計画のステップであり、結合の図ではノードとして表すことができます。実際の実行は、いくつかのパラメータ設定および結合される小さな各表のカーディナリティによってわずかに異なる場合があります。

簡単な 2 つの表の結合の詳細

特定の 2 つの表に対する結合操作の場合、一方の表が大きい表と判断され、他方が小さい表と判断されます。この判断は、表内のブロック数および述語のカーディナリティの推定値に基づきます。この結合は、2 つのバッチプリミティブステップ (BPS) を使用して実行されます。小さい表をスキャンする BPS は、UM にアクセスしてデータを戻すときに、利用可能な任意のフィルタを適用します。これは、BPS (小) と表すことができます。

次の条件付き動作は、大規模な結合の実行に必要なネットワーク通信またはプロセス間通信を最小化する（データの送信コストを最小限に抑える）ために行われます。データが UM に戻されると、データセットのサイズが判断されます。この測定値は、次のパラメータと比較されます。

PmMaxMemorySmallSide

デフォルト値は 64M (サイズは MB 単位) です。

小さい表の 100 万程度の行が、大きい表の数十億 (あるいは数兆) の行に対して結合される場合には、このデフォルトで十分です。実際のカーディナリティは、結合データ型と、問合せに含まれる小さい表の追加列によって異なります。データセットのサイズが PmMaxMemorySmallSide よりも小さい場合は、データが PM に送信されて分散ハッシュマップが作成されます。そうではない場合、UM で作成されます。

UM または PM で小さい表のハッシュマップがインスタンス化されると、最大の表に対して BPS (大) が発行されます。小さい表のハッシュマップが PM に送信された場合、結合操作および集計操作は完全分散で行われます。小さい表のハッシュマップが UM にある場合、結合および集計は UM で行われます。どの場合も、列および表のフィルタは、PM 上で、完全にマルチスレッド化および分散された方法で同様に適用されます。

実質的に、2 つの表の結合は、2 つのバッチプリミティブステップによって実行されます。

BPS (小) → BPS (大)

データウェアハウスの用語を使用すると、この簡単な問合せは次のように言い換えられます。

BPS (ディメンション) → BPS (ファクト)

複数結合の詳細

単一の BPS で複数の小さい表と 1 つの大きい表を結合する場合にも、この結合をサポートするために前述の処理モデルが使用されます。

1 つの大きな表が 2 つの小さい表と結合される場合、操作は次のようになります。

BPS (小_1)

\

BPS (小_2) ----> BPS (大)

これは、大きい表を処理する前に、小さい表の 2 つのステップで実行されます。この処理モデルは、小さい表のハッシュマップが UM または PM のどちらでインスタンス化されたかにかかわらず、同じであることに注意してください。実際に、前述の図のステップは、各オブジェクトのサイズにかかわらず実行できます。実行時の各オブジェクトの正確なカーディナリティやサイズに応じて、両方の結合を PM または UM で実行したり、任意の組合せに分割して実行することができます。

実際、この結合処理モデルは、任意の数の結合操作を処理できます。この際、結合操作に含まれる小さい各表に対して 1 つの BPS が使用され、最大の表に対して 1 つの BPS が使用されます。

スタースキーマのデータモデルの場合は、この方法で、自己結合を使用しない多数の問合せに対応することができます。

```
BPS (ディメンション_1)    \  
BPS (ディメンション_2)    \  
BPS (ディメンション_3)    > BPS (ファクト)  
  中略                      /  
BPS (ディメンション_20)   /
```

これによって、スレッド間の同期を最小限に抑えた状態で、完全にマルチスレッド化および分散（オプション）された処理を実行できます。また、ファクト表から行または集計を戻す処理は、可能なすべてのフィルタが適用されるまで保留されます。

その他の結合の組合せは、連鎖結合演算子および複数結合演算子の様々な組合せを使用して行われます。たとえば、スノーフレイクスキーマでは、BPS（ファクト）表用のハッシュマップを構築するための前提条件である追加の結合操作が発生します。内部結合操作および外部結合操作の両方がこれらの連鎖結合演算子および複数結合演算子でサポートされています。

結合の最適化

InfiniDB オプティマイザは、統計を使用してグローバル結合ツリー内で最大の表を判断し、その表をストリーム表として使用できるように操作を順序付けます（データの転送コストを最小限に抑えます）。グローバルで最大の表によってまだ設定されていない追加のサブツリーでも、大きい表または小さい表のどちらであるかをローカルで判断できます。実際のカーディナリティが類似している場合、表の選択はパフォーマンスに影響しないことがあります。大きい1つのファクト表と小さい複数のディメンション表を含むデータセットの場合、準最適なものを選択される可能性はかなり小さくなります。オプティマイザによって行われた判断が準最適である場合に問合せをチューニングできるようにするには、結合操作内で最大の表を設定するためのヒントを利用できます。

INFINIDB_ORDERED ヒントには、From 句の最初の表はグローバルで最大の表として処理され、可能な場合はその表の結合結果のマテリアライズが除外されるように結合が順序付けられることが示されます。

このヒントは、グローバルで最大の表を結合対象の最後の表として設定します。問合せ内でこのグローバルで最大の表に結合された表により、その最後の操作に含める表（または結合サブツリー）が決定されます（From 句内の表の順序ではありません）。InfiniDB では、複数の結合操作が1つのBPSで実行されます。そのため、1つの大きい表とn個の小さい表の結合では、ほとんどの場合、小さい表の結合の順序によって問合せのI/O特性が変わることはありません。

たとえば、このヒントを次のように使用すると、オプティマイザは実際のカーディナリティに関係なく、region 表を結合ツリー内で最大の表として処理します。

```
select /*!INFINIDB_ORDERED */ r_regionkey
  from region r, customer c, nation n
 where r.r_regionkey = n.n_regionkey
       and n.n_nationkey = c.c_nationkey;
```

ネステッドループ操作は行われなかったため、この場合も、ヒントによって駆動表が指定されたり、1 つずつ結合する表の順序が決定されることはありません。これらの概念は、少なくとも従来の使用方法では InfiniDB 内に存在しません。

主要なチューニングパラメータ : PmMaxMemorySmallSide

PmMaxMemorySmallSide パラメータによって、PM に送信される小さい表の単一ハッシュマップのサイズに対する上限が設定され、結合を完全分散で実行するかどうかが決まります (エンタープライズ版の場合)。

単一サーバーのインストールの一般的なチューニングガイドライン

PmMaxMemorySmallSide は、小さい表の結合について予測できる最大サイズをサポートできる大きさに設定します。ただし、利用可能なメモリーを上限とします。単一サーバー内でも、データの送信コストを考慮する必要があり、先に (PM で) フィルタを実行すると、サーバー内部のデータ送信コストが削減されます。

複数 PM のインストールの一般的なチューニングガイドライン

複数 PM 構成の PmMaxMemorySmallSide のチューニングは、サーバーで利用可能なメモリーの量、データバッファキャッシュ用に必要なメモリーの量、同時に発生することが予測される同時結合操作の数、およびこれらの結合操作のサイズに関係します。

次のようにチューニングすることをお勧めします。

(同時実行される PM の小さい表のハッシュマップの数) * (各 PM のハッシュマップの平均サイズ)

前述の値が次より小さくなる必要があります

(サーバーのメモリーの合計) - (データバッファキャッシュのサイズ)

このパラメータを設定することによってメモリーが明示的に消費されるわけではないことに注意してください。このパラメータでは、PM でインスタンス化できるハッシュマップの最大サイズのみが制限されます。

16GB のメモリーを備えたサーバーで、同時実行はそれほど多くなく、2 番目に大きい表の結合カーディナリティが最大で 1000 万である場合は、PmMaxMemorySmallSide を 512M に設定して、データバッファキャッシュを 14GB に設定すると適切です。

同時実行がかなり多い場合でも、各 PM のハッシュマップの平均サイズによっては、512M の値で有効な場合があります。ただし、サーバーのメモリー利用率が 100% に近づいた場合、またはこの値を越えてスワップが開始された場合は、PmMaxMemorySmallSide の設定を小さくすると、サーバーのメモリーが削減される可能性があります。

メモリー管理

主要なチューニングパラメータ : NumBlocksPct および TotalUmMemory

各パフォーマンスモジュール処理内のデータバッファキャッシュ専用のメモリーの量は、次の Calpont.xml パラメータによって設定します。

PM のメモリーパラメータ :

NumBlocksPct

デフォルトは提供形態によって異なります (単一サーバーのデフォルトは 50 で、複数サーバーのインストールのデフォルトは 86 になります)。値が 86 の場合、データバッファキャッシュは、サーバーで利用可能なメモリーの 86% まで消費できます。

ユーザーモジュールとパフォーマンスモジュールが同じサーバーで実行される単一サーバー構成の場合は、NumBlocksPct の設定を低くする必要があります。これは、ユーザーモジュールが、一時データセットを管理するために一時領域を著しく必要とする場合があるためです。

TotalUmMemory パラメータは、ユーザーモジュールでの結合、集計および集合操作の中間結果の管理に利用可能なメモリーの最大量を制限します。これは、メモリーを専用に確保するのではなく、単一結合のハッシュマップのサポートで消費可能なメモリーの最大量を制限します。通常は、必要に応じてこの設定を引き上げることで、何億もの行が含まれることのある小さい表の臨時ハッシュマップを最小限の影響で使用できます。

PM 処理と UM 処理が異なるサーバーに分離されるエンタープライズ版インストールでは、小さい表の平均の結合カーディナリティが 100K より小さい場合、キャッシュを PM の合計メモリーの 95% に設定しても差し支えありません。この場合、PM のメモリー利用率は、様々な同時実行のシナリオに対して 95-97% の間で変動します。

InfiniDB では、大きい表の数十億または数兆の行を小さい表の任意の数のハッシュマップに結合することがサポートされており、大きい表のサイズは、どのパラメータ設定でもまったく制限されないことに注意してください。TotalUmMemory パラメータは小さい表のマップにのみ適用されます。

スケーラビリティ

スケーラビリティ：データのサイズとパフォーマンス

特定のシステムが構成され、問題となる問合せが最適化されて効率的に実行されたら、パフォーマンスモジュールの規模を調整することによってパフォーマンスを向上することができます。また、システムの規模を調整することで、大容量のデータを一貫したパフォーマンスで分析できる場合があります。

スケーラビリティ：ユーザーモジュール

ユーザーモジュール上で CPU 使用率が増加している場合、その負荷を分散するためにユーザーモジュールを追加できます。InfiniDB の目標の 1 つはできるかぎり多くの作業を分散することであるため、多くの問合せが PM 上で 99% の CPU サイクルを実現できます。

ツールおよびユーティリティのチューニング

問合せのサマリー統計 : calgetstats

以前の操作のサマリー結果を表示する問合せの後に、`select calgetstats()` コマンドを実行できます。次の問合せの例では、3つのフィルタ（そのうち2つは結合で表されている）を使用して、ファクト表に 59 億 9000 万行が含まれているスタースキーマベンチマークのデータセットに対して4つの表の結合が実行されています。最後の集計は約 2000 万行に対して行われ、この例では 8PM 構成が使用されています。

結果の詳細は、出力例のとおりです。

```
select  d_year, lo_tax, p_size, s_region, count(*)
from    dateinfo, part, supplier, lineorder
where   s_suppkey = lo_suppkey
        and d_datekey = lo_orderdate
        and p_partkey = lo_partkey
        and lo_orderdate between 19980101 and 19981231
        and s_nation = 'BRAZIL'
        and p_size <> 23
group by 1,2,3,4
order by 1,2,3,4;
```

```
+-----+-----+-----+-----+-----+
| d_year | lo_tax | p_size | s_region | count(*) |
+-----+-----+-----+-----+-----+
| 1998   | 0.00  | 1      | AMERICA  | 47607   |
| 1998   | 0.00  | 2      | AMERICA  | 47194   |
... results abbreviated ...
| 1998   | 8.00  | 48     | AMERICA  | 47846   |
| 1998   | 8.00  | 49     | AMERICA  | 47394   |
| 1998   | 8.00  | 50     | AMERICA  | 47030   |
+-----+-----+-----+-----+-----+
441 rows in set (3.66 sec)
```

```
mysql> select calgetstats();
```

```
+-----+
-----+
-----+
-----+
-----+
```

```

| calgettats()
|
+-----+
+-----+
| Query Stats: MaxMemPct-4; NumTempFiles-0; TempFileSpace-0MB; ApproxPhyI/O-0;
CacheI/O-1363254; BlocksTouched-1363254; PartitionBlocksEliminated-2637824;
MsgBytesIn-811MB; MsgBytesOut-0MB; Mode-Distributed| 1256555629 961957 |
+-----+
+-----+
1 row in set (0.02 sec)

```

calgetstats() 関数の出力は次のとおりです。

- ◆ **MaxMemPct-4** : このフィールドには、ユーザーモジュール (UM) の結合、グループ化、集計、重複行の削除などの操作を行うために使用されている UM のメモリー利用率が示されます。
- ◆ **NumTempFiles-0** : このフィールドには、ユーザーモジュール (UM) の結合、グループ化、集計、重複行の削除などの操作を行うために使用されている UM の一時ファイル使用数が示されます。
- ◆ **TempFileSpace-0MB** : このフィールドには、ユーザーモジュール (UM) の結合、グループ化、集計、重複行の削除などの操作を行うために使用されている UM の一時ファイル使用サイズが示されます。
- ◆ **ApproxPhyI/O-0** : このフィールドには、問合せに対するストレージからの読取りが示されます。状況によっては、実際とは少し異なる場合があります (通常は 0.1% 未満)。
- ◆ **CacheI/O-1363254** : このフィールドには、リクエストされた個別の物理 I/O 読取りの数が削減されている、問合せで必要なブロックアクセスが示されます。
- ◆ **BlocksTouched-1363254** : このフィールドには、問合せで必要なブロックアクセスが示されます。
- ◆ **PartitionBlocksEliminated-2637824** : このフィールドには、エクステンツマップの最小値または最大値に基づいて発生したパーティションの除外によって排除されたブロックが示されます。フィルタが適用されるまで、列ストレージまたは遅延している列の読取りでの I/O の削減はレポートされないことに注意してください。
- ◆ **MsgBytesIn-811MB** : データ移動を処理するプロセスの測定単位。
- ◆ **MsgBytesOut-0MB** : データ移動を処理するプロセスの測定単位。
- ◆ **Mode-Distributed** : 問合せ結合処理が InfiniDB 内で処理されたか、従来の MySQL の結合機能で処理されたかを示すインジケータ。大規模なデータに対して最高のパフォーマンスを実現するには、この処理を分散することをお勧めします。

問合せの詳細統計 : calgettrace

InfiniDB の SQL トレース機能を有効にすることでさらに詳細なレベルでの追加情報を表示できます。SQL トレースの使用手順は次のとおりです。

1. コマンド `select calsettrace(1);` を発行して、新しいトレースを有効にします。
2. 問合せを実行します。
3. コマンド `select calgettrace();` を発行して、トレースの結果を表示します。

たとえば、新しいトレースを有効にした後、次のような問合せを実行して分析できます。

```
select d_year, lo_tax, p_size, s_region, count(*)
  from dateinfo, part, supplier, lineorder
 where s_suppkey = lo_suppkey
       and d_datekey = lo_orderdate
       and p_partkey = lo_partkey
       and lo_orderdate between 19980101 and 19981231
       and s_nation = 'BRAZIL'
       and p_size <> 23
 group by 1,2,3,4
 order by 1,2,3,4;
```

```
mysql> select calgettrace();
+-----+
| calgettrace()
|
+-----+
|
Desc Mode Table      TableOID ReferencedOIDs      PIO    LIO    PBE    Elapsed Rows
BPS  PM   part        3513    (3521,3514)         1370   1372   0      0.281  1371884
BPS  PM   dateinfo    3558    (3559,3563)         12     10     0      0.126  2556
DSS  PM   supplier    3528    (3539)              0      1      -      0.000  1
BPS  PM   supplier    3528    (3533,3529,3540)   2449   3551   0      0.263  40078
HJS  PM   supplier    3528    -                   -      -      -      ----- -
      -supplier
BPS  PM   lineorder   3493    (3499,3497,3498,3508) 132060 131298 266240 5.553  1489347
HJS  PM   lineorder   3493    -                   -      -      -      ----- -
      -dateinfo
HJS  PM   lineorder   3493    -                   -      -      -      ----- -
      -part
HJS  PM   lineorder   3493    -                   -      -      -      ----- -
      -supplier
TAS  UM   -           -       -                   -      -      -      5.374  441
TAS  UM   -           -       -                   -      -      -      5.374  441
+-----+
1 row in set (0.02 sec)
```

calgettrace () の出力には次のヘッダーが含まれています。

- ◆ Desc : 実行された操作。
- ◆ Mode : UM 内または PM 内のいずれで実行されたか。
- ◆ Table : 列をスキャンまたは投影できる表。
- ◆ TableOID : スキャンされた表のオブジェクト ID。
- ◆ ReferencedOIDs : 問合せで必要な列のオブジェクト ID。
- ◆ PIO : 問合せで実行された物理 I/O (ストレージからの読取り)。
- ◆ LIO : 問合せで実行された論理 I/O。ブロックアクセスとも呼ばれます。
- ◆ PBE : 除外されたパーティションブロック (PBE) では、エクステンタマップの最小値または最大値によって除外されたブロックが識別されます。
- ◆ Elapsed : 表示されているステップの経過時間。
- ◆ Rows : 返された中間行

トレース出力には、操作の順序の他、ブロックアクセス、経過時間および行カーディナリティの操作コストがレポートされます。最初のフィルタは supplier 表に対して適用されます。

```
◆ DSS PM supplier 3528 (3539) 0 1 - 0.000 1
```

= 'Brazil' のフィルタは、複数の行で 1 つの文字列値を共有できるディクショナリ構造内に格納されている可変長フィールドを参照します。これは、ディクショナリシグネチャステップ (DSS) という個別のステップで行われます。

```
◆ BPS PM supplier 3528 (3533,3529,3540) 2449 3551 0 0.263 40078
```

この BPS は、後続の結合で supplier から他の列を読み取ることを示します。

```
◆ HJS PM lineorder 3493 - - - - -  
-supplier
```

このハッシュ結合ステップ (HJS) は、supplier の投影と DSS フィルタを関連付けます。これは、実際には前の BPS の一部であり、時間には常に 0 (ゼロ) がレポートされることに注意してください。Mode フィールドは PM 結合がここで実行されたことを示しています。

```
◆ BPS PM dateinfo 3558 (3559,3563) 12 10 0 0.126 2556
```

この BPS は、後続の結合で dateinfo 表の列を読み取ることを示します。

```
◆ BPS PM part 3513 (3521,3514) 1370 1372 0 0.281 1371884
```

この BPS は、後続の結合で part 表の列を読み取ること示します。

```
◆ BPS PM lineorder 3493 (3499,3497,3498,3508) 132060 131298 266240 5.553
1489347
```

この BPS は、lineorder からスキャンおよび投影を行い、supplier、dateinfo および part に対して結合操作を実行することを示します。すべてのフィルタおよび結合が終了した後の出力カーディナリティは 21,340,320 行です。

```
◆ HJS PM lineorder 3493 - - - - -
-dateinfo
```

HJS は、前の BPS 内で発生するハッシュ結合操作のインジケータです。Mode フィールドの値 PM または UM は、UM または PM にハッシュマップが作成されたかどうかを示します。ハッシュ結合操作は実際には前の BPS 内で発生するため、ここにレポートされる経過時間は常に 0 (ゼロ) になります。

```
◆ TAS UM - - - - - 5.374 441
```

タプル集計ステップ (TAS) では、UM が入力 BPS から最初の中間集計結果を受信してから必要な集計を完了するまでの経過時間がレポートされます。ほとんどの集計シナリオで、この時間は前の BPS 操作と密接に対応します。この操作では、レポートされる Mode は常に UM となりますが、基本機能は常に 2 フェーズ操作 (まず PM でローカル集計が実行され、UM で最終集計が実行される) で行われます。

多くの操作が並行して行われるため、個々の経過時間の合計が問合せの経過時間よりも長くなることに注意してください。前述の例では、BPS ステップに 5.553 秒、タプル集計ステップ (TAS) に 5.374 秒かかっていますが、実際にはパイプライン操作として行われており、BPS は部分的に集計されたデータを TAS に入力します。したがって、BPS (lineorder) とその後の TAS の両方を実行する場合の経過時間はレポートされる 5.553 秒により近くなります。

キャッシュのフラッシュ : calflushcache

InfiniDB では、すべての PM 間のデータバッファキャッシュからデータをフラッシュする物理 I/O の簡単なテストを実行できる開発ユーティリティまたはテストユーティリティが提供されています。これは、次のコマンドを発行することで実行されます。

```
select calflushcache();
```

これは、開発ユーティリティまたはテストユーティリティとして使用することのみを目的としています。キャッシュのフラッシュおよび PIO の増加によるメリットは報告されていません。それどころか、通常は問合せの実行時間が長くなります。

パラメータの読取りまたは変更 : configxml.sh

getconfig または setconfig によって Calpont.xml パラメータファイル内の値の読取りまたは設定を行うことができるユーティリティが提供されています。

使用方法 : /usr/local/Calpont/bin/configxml.sh {setconfig|getconfig} section
variable set-value

このマニュアル内で説明されているパラメータを更新する構文の例 :

```
/usr/local/Calpont/bin/configxml.sh setconfig HashJoin PmMaxMemorySmallSide 640M  
/usr/local/Calpont/bin/configxml.sh setconfig JobList MaxOutstandingRequests 7  
/usr/local/Calpont/bin/configxml.sh setconfig DBBC NumBlocksPct 90
```

エクステントマップの表示 : editem

editem というユーティリティを使用すると、指定した列にデータが移入されているかどうかを判断できます。このユーティリティではオブジェクト ID が必要です。オブジェクト ID は、トレース出力またはシステムカタログのいずれかから取得できます。

警告 : editem ユーティリティは、検査用の値の表示には安全に使用できます。データが破損する恐れがあるため、editem を他の目的に使用する場合は、Calpont サポートの指示に従って使用するか、隔離されたテスト環境でのみ使用してください。

システムカタログからの選択の例 :

```
select columnname, objectid from calpontsys.syscolumn where  
tablename = 'lineorder' and columnname = 'lo_commitdate';
```

editem のヘルプでは、提供されている他のすべてのユーティリティと同様に、-h フラグを使用することで (例 : /usr/local/Calpont/bin/editem -h)、使用可能なコマンドのリストが提供されます。

editem -o <objectid> コマンドによって、指定した列の最小値および最大値が表示されます。最小値が存在しない場合はそのデータ型に対して可能な最大値で表され、最大値が存在しない場合はそのデータ型に対して可能な最小値で表されることに注意してください。

列データストレージの相違点

列ストレージのメリットとトレードオフを次に示します。トレードオフを回避したり、列ストレージ機能を最大限に活用する方法についても説明します。

- ◆ **挿入のトレードオフ**：列ストレージでは、個々の行を挿入する際のパフォーマンスプロファイルが異なります。任意の行ベースのシステムに 1 行挿入する場合、アクセスする可能性があるのは 1 つの表ブロックおよび複数の索引ブロックです。列ストレージの DBMS に挿入する場合は、列ごとに 1 つ以上のブロックにアクセスしますが、索引ブロックに対する追加コストがありません。このコストは、バルクロード (cpimport)、LOAD DATA INFILE またはバルク挿入操作では変わることにご注意ください。何千行も挿入する場合、関連するブロックアクセスは数千行に集中し、行 DBMS の表および索引の組合せと比較すると、列ストレージは数千を超えるとより効率的になります。
 - **InfiniDB の方法**: cpimport を使用します。InfiniDB の cpimport バルクロードユーティリティを使用すると、列ストレージ (索引を使用しない) での挿入に必要なブロックは、行ベースの DBMS の場合より少なくなります。また、索引を使用しないと、通常、より高速のロードパフォーマンスを実現でき、大規模なデータのロード時間の一貫性が向上します。InfiniDB では、表が空でも、表に 100 億行が含まれていても、一連の行のロードは同じ操作です。
- ◆ **削除のトレードオフ**：列ストレージでは、個々の行の削除のパフォーマンスプロファイルが異なります。任意の行ベースのシステムから 1 行削除する場合、アクセスする可能性があるのは 1 つの表ブロックおよび複数の索引ブロックです。列ストレージの DBMS から削除する場合は、列ごとに 1 つ以上のブロックにアクセスしますが、索引ブロックに対する追加コストがありません。特定の範囲の行から数千行を削除する場合は、関連するブロックアクセスが集中し、行 DBMS の表および索引の組合せと比較すると、列ストレージはより効率的になります。

削除操作のパフォーマンスは、どの DBMS でも、削除する行がブロック間で分散されている状態によって異なる場合があります。たとえば、8K ブロックごとに 100 行を格納する行ベースの DBMS から 1000 行を削除する場合は、10 個 (最適な分散) から 1000 個 (最悪の分散) の範囲のブロックにアクセスする可能性があります。列ストレージでは、10 列の表 (各列がブロックごとに 1024-8192 行を格納する) の場合、最良のシナリオは 10 個のブロック (最適な分散) ですが、最悪の場合は 10,000 個のブロック (最悪のシナリオ) です。

- **InfiniDB の方法** : バッチで削除します。たとえば、100 万行を指す 100 個のキーを含む IN 句に基づいた行の削除は、100 個の個別の削除文より大幅に短時間で実行されます。
- ◆ **更新のメリット** : 列ストレージの更新では、ほとんどのシナリオで大きなメリットが得られます。1000 行の行 DBMS の更新の最良のシナリオは、100 個のブロック（列ごとに 100 行と想定）です。列 DBMS の最良のシナリオは 1 個のブロックです。索引も更新する必要がある場合は、パフォーマンス上のメリットが増えます。
- **InfiniDB の方法** : すでに最適化されているため、通常は更新が非常に高速に実行されます。

データのロード速度およびリアルタイムに近いロード

データのロード速度は、バッチ処理の状況、基礎となるストレージ機能、表定義、データ型および値によって異なります。ただし、一般的な目安として、cpimport バルクロードユーティリティのロード速度は、1 秒当たり数十万から数百万行になる場合があります。LOAD DATA INFILE のロード速度は、1 秒当たり数千行の場合があります。個々の挿入のロード速度はさらに遅く、1 秒当たり数十行です。どの場合も、データは読取り一貫性動作が確保されるモデルで使用可能になります。

cpimport バルクロード機能を使用すると、一度に 1000 行、数百万あるいはそれ以上のデータをロードできます。一度に 10 万行以上をロードするときに速度がピークになるインポートでは、スタートアップ時間に約 1 秒程度かかります。基本的な操作は変わらないため、このロード速度は以降のロードでは一貫しています。

どの方法でデータをロードする場合も、行の挿入によってエクステントが表に追加される可能性があります。この際、エクステントがストレージ内で全体に連続して配置されるようにするため、通常は 5-20 秒の遅延が発生します。これによって、マルチブロック読取りのメリットを最大化し、以降の選択操作で最大限のパフォーマンスが確保されます。

処理の優先順位付け

デフォルトでは、インポート処理と問合せ処理の優先順位は同じです。この優先順位は、InfiniDB 構成ファイル Calpont.XML にエントリを追加することによって変更できます (Linux のみ)。これは、同時実行の問合せとインポートが定期的に行われ、結果としてインポート速度が低下する場合に役立つことがあります。このエントリは次のとおりです。

```
<ExeMgr1>
  <Priority>###</Priority>
</ExeMgr1>

<PrimitiveServers>
  <Priority>###</Priority>
</PrimitiveServers>

<WriteEngine>
  <Priority>###</Priority>
</WriteEngine>
```

の値は次のとおりです。

- ◆ 1-40 : 1 は最も低い優先順位で、40 は最も高い優先順位です。前述の 3 つのプロセスのデフォルトの優先順位は 21 です。

ExeMgr1 および PrimitiveServers の優先順位を変更すると問合せの優先順位が変更され、WriteEngine の優先順位を変更するとインポートの優先順位が変更されます。

InfiniDB のパフォーマンスの目安

熟練したチューニング担当者が行ベースの DBMS システムでの経験に基づいて適切と感じる多数の目安があります。InfiniDB ではこの従来の目安のいくつかは大幅に変わっており、実際、より適切にデータウェアハウスの複雑さに対処できる多数の新しい方法を示している場合があります。

- ◆ 5% 未満のデータを問い合わせるときには索引が役立ちます。この目安は、表をスキャンしない（全表スキャンを行わない）列ストレージでは変わります。たとえば、類似する 20 列を含む表の場合、1 列をスキャンするコストが表をスキャンするコストの 5% になるため、目安は 0.25% 未満のデータとなります。索引の有用性の低下、ランダムな I/O 動作および多大な維持コストを理由として、現時点では InfiniDB で索引は実装されていません。
 - InfiniDB の新しいパラダイム：I/O がより効率的であるマルチスレッド操作や分散操作による処理が利用されます。

- ◆ 大きい表をモデリングするときは、アクセス頻度が低い列を含めないでください。行ベースの DBMS システムの場合、表の一部として作成される追加の列によって、多数の行を含む（また、Covering Index またはデータのその他の複製やマテリアライズで処理できない）**すべての**問い合わせの速度が低下する可能性があります。
 - InfiniDB の新しいパラダイム：問い合わせで参照されない列は無視されるため、追加のストレージおよび負荷に対するコストのみで、不定期の分析に追加のデータを使用できます。

- ◆ 列のデータ型はそれほど重要ではありません。多くの行ベースの DBMS システム（特に、すべての列を可変サイズのデータ型に格納するシステム）では、データ型に基づくパフォーマンスの違いがない場合があります。データに単一の文字値のみが含まれる場合は、行ベースの DBMS の char(2) を char(1) に変更すると、索引なしの列をスキャンするコストが変わる可能性があります、その程度は 1% 以下です。列ストレージシステムの場合、列のスキャンに必要なブロックは半分であるため、効率が向上して列ごとの検索が高速化されます。
 - InfiniDB の新しいパラダイム：索引をチューニングするのではなく、列のデータ型をチューニングして検索を高速化できます。

- ◆ 冗長データに利点はありません。冗長データによって、表にアクセスするすべての問合せのコストが増えるだけでなく、更新または一貫性のロジックに関連する問題が必ず発生します。ただし、データのロードに一度だけ書込み可能な方式を使用して実装されるデータウェアハウスのシナリオでは、更新または一貫性の問題は発生しない場合があります。行ベースのストレージでは、表にアクセスするすべての問合せに対して追加の I/O コストが依然として必要となります。
 - InfiniDB の新しいパラダイム：追加の列によってデータへの新しいアクセスパスが提供される場合があります（追加の列がフィールドの主要な部分か、複数のフィールドの連結であるかは関係ありません）。

- ◆ 有用なキャッシュにするには、キャッシュをファクト表（またはアクティブなパーティション）より大きくする必要があります。行ベースの独自の DBMS のほとんどでは、2 回目のスキャンがキャッシュによって処理されるようにするには、全表スキャンまたは表内のアクティブなパーティションの全体スキャンがメモリーに完全に収まる必要があります。
 - InfiniDB の新しいパラダイム：列ストレージでは、問合せで参照されない列が削除され、列へのアクセスが個別に行われるため、キャッシュのサイズが表のサイズのほんの一部にすぎない場合でもキャッシュのメリットが得られます。

- ◆ ディスクからの読取りは、キャッシュからの読取りに比べて 20-40 倍低速です。実際の比率は多数の要因によって異なりますが、通常、ディスクから読み取るときはパフォーマンスが大幅に低下します。20-40 倍の比率は、問合せコストの 95%-97.5% がディスクからの読取りに直接関連することを示します。
 - InfiniDB の新しいパラダイム：キャッシュにより 2-10 倍改善します。InfiniDB の I/O の効率性（列ストレージ、パーティションブロックの除外、フィルタ後までの I/O の遅延）により、ほとんどの問合せの絶対 I/O コストが大幅に減ります。結果として、キャッシュにより相対コストが 2-10 倍速くなります。

- ◆ パーティション化はデータウェアハウスのパフォーマンスにとって重要です。これは、実際に今も当てはまります。従来の独自の DBMS システムでは、DBA は 1 つまたは 2 つのレベルのパーティション化方法を宣言して、この 1 または 2 列に対して適切なフィルタを使用して問合せのブロックアクセスを減らすことができます。
 - InfiniDB の新しいパラダイム：列ストレージに基づいて、自動パーティション化を使用できます。ロード方法に基づいて、データの昇順化やその他のクラスタ化が行われているすべての列に対して自動パーティション化を使用できます。これは 1 または 2 列に限定されず、適切な特性を持つすべての列に使用できます。
 - InfiniDB パーティションは必ずしも実行する必要はありません。データがランダムにロードされる場合は、パーティションをなくすことによるメリットはない可能性があります。

- ◆ 結合には適切な索引が必要です。どのネステッドループ結合操作でも、ループ内の検索が索引によってサポートされると、パフォーマンスが向上します。多くの場合、索引を使用しないと結合のパフォーマンスは大幅に低下します。また、ネステッドループ操作を反転した場合は、別の索引が必要になります。
 - InfiniDB の新しいパラダイム：ハッシュ結合機能を利用して索引の必要性をなくし、ネステッドループ処理に関連する大量の行単位処理のコストをなくします。

- ◆ 条件付きのパフォーマンス期待値。適切にチューニングし、チューニングによって定義されている境界内に問合せが存在する場合は、行ベースの DBMS システムをこれらの問合せ用にチューニングできます。ただし、索引を使用しない、または明示的に宣言された 1 つまたは 2 つのパーティション列によって解決されないその他の問合せのパフォーマンスは非常に低い場合があります、パフォーマンスの低下が 100 倍になることもあります。索引が完全にキャッシュされる場合と、ランダムなアクセスを処理するために何度も再ロードされる場合とでは、索引のパフォーマンス特性が大きく異なります。結合の順序によってネステッドループ操作のコストが大幅に変わる可能性があります。
 - InfiniDB の新しいパラダイム：より一貫したパフォーマンスが実現します。10 億のレコードに対して列が 1、2、4 または 8 バイトの表の場合はスキャン率にいくらかの違いがありますが、桁違いになることはありません。結合の順序は、データの送信コストを最小限にするように自動的に処理されます。

- ◆ すべての列が必要なわけではない場合、「Select *」に利点はありません。
 - InfiniDB の従来どおりのパラダイム：すべての列が必要なわけではない場合、従来どおり利点はありません。

追加リソース、ダウンロードおよびサポート

オープンソース版の InfiniDB は、コミュニティ版の Web サイト (www.infinidb.org) からダウンロードできます。製品版の InfiniDB については、infinidb_doc@ashisuto.co.jp にご連絡ください。

GNU Free Documentation License

GNU Free Documentation License
Version 1.3, 3 November 2008

Copyright © 2000, 2001, 2002, 2007, 2008 Free Software Foundation, Inc. <<http://fsf.org/>>

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

0. PREAMBLE

The purpose of this License is to make a manual, textbook, or other functional and useful document "free" in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondly, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of "copyleft", which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

1. APPLICABILITY AND DEFINITIONS

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The "Document", below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as "you". You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A "Modified Version" of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A "Secondary Section" is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document's overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The "Invariant Sections" are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The "Cover Texts" are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

A "Transparent" copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not "Transparent" is called "Opaque".

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The "Title Page" means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, "Title Page" means the text near the most prominent appearance of the work's title, preceding the beginning of the body of the text.

The "publisher" means any person or entity that distributes copies of the Document to the public.

A section "Entitled XYZ" means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as "Acknowledgements", "Dedications", "Endorsements", or "History".) To "Preserve the Title" of such a section when you modify the Document means that it remains a section "Entitled XYZ" according to this definition.

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

2. VERBATIM COPYING

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

3. COPYING IN QUANTITY

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each

Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

4. MODIFICATIONS

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

- ◆ A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
- ◆ B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.
- ◆ C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
- ◆ D. Preserve all the copyright notices of the Document.
- ◆ E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
- ◆ F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
- ◆ G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
- ◆ H. Include an unaltered copy of this License.
- ◆ I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document

as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.

- ◆ J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
- ◆ K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.
- ◆ L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
- ◆ M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.
- ◆ N. Do not retitle any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.
- ◆ O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties—for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

5. COMBINING DOCUMENTS

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled "History" in the various original documents, forming one section Entitled "History"; likewise combine any sections Entitled "Acknowledgements", and any sections Entitled "Dedications". You must delete all sections Entitled "Endorsements".

6. COLLECTIONS OF DOCUMENTS

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

7. AGGREGATION WITH INDEPENDENT WORKS

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an "aggregate" if the copyright resulting from the compilation is not used to limit the legal rights of the compilation's users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document's Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

8. TRANSLATION

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations

requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled "Acknowledgements", "Dedications", or "History", the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title.

9. TERMINATION

You may not copy, modify, sublicense, or distribute the Document except as expressly provided under this License. Any attempt otherwise to copy, modify, sublicense, or distribute it is void, and will automatically terminate your rights under this License.

However, if you cease all violation of this License, then your license from a particular copyright holder is reinstated (a) provisionally, unless and until the copyright holder explicitly and finally terminates your license, and (b) permanently, if the copyright holder fails to notify you of the violation by some reasonable means prior to 60 days after the cessation.

Moreover, your license from a particular copyright holder is reinstated permanently if the copyright holder notifies you of the violation by some reasonable means, this is the first time you have received notice of violation of this License (for any work) from that copyright holder, and you cure the violation prior to 30 days after your receipt of the notice.

Termination of your rights under this section does not terminate the licenses of parties who have received copies or rights from you under this License. If your rights have been terminated and not permanently reinstated, receipt of a copy of some or all of the same material does not give you any rights to use it.

10. FUTURE REVISIONS OF THIS LICENSE

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <http://www.gnu.org/copyleft/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License "or any later version" applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free

Software Foundation. If the Document specifies that a proxy can decide which future versions of this License can be used, that proxy's public statement of acceptance of a version permanently authorizes you to choose that version for the Document.

11. RELICENSING

"Massive Multiauthor Collaboration Site" (or "MMC Site") means any World Wide Web server that publishes copyrightable works and also provides prominent facilities for anybody to edit those works. A public wiki that anybody can edit is an example of such a server. A "Massive Multiauthor Collaboration" (or "MMC") contained in the site means any set of copyrightable works thus published on the MMC site.

"CC-BY-SA" means the Creative Commons Attribution-Share Alike 3.0 license published by Creative Commons Corporation, a not-for-profit corporation with a principal place of business in San Francisco, California, as well as future copyleft versions of that license published by that same organization.

"Incorporate" means to publish or republish a Document, in whole or in part, as part of another Document.

An MMC is "eligible for relicensing" if it is licensed under this License, and if all works that were first published under this License somewhere other than this MMC, and subsequently incorporated in whole or in part into the MMC, (1) had no cover texts or invariant sections, and (2) were thus incorporated prior to November 1, 2008.

The operator of an MMC Site may republish an MMC contained in the site under CC-BY-SA on the same site at any time before August 1, 2009, provided the MMC is eligible for relicensing.