

Snowflake技術仕様書

序文

本技術仕様書は、クラウドデータウェアハウスソリューションであるSnowflakeについて、そのアーキテクチャ、主要機能、およびベストプラクティスを詳細に記述することを目的としています。Snowflakeは、その柔軟性、スケーラビリティ、およびパフォーマンスにより、現代のデータ駆動型ビジネスにおいて不可欠なツールとなっています。本ドキュメントは、Snowflakeの導入を検討している組織や、既存のSnowflake環境を最適化したいと考えているユーザーに対し、包括的な情報を提供します。

Snowflakeは、従来のデータウェアハウスとは異なり、ストレージとコンピュー트를分離した独自のアーキテクチャを採用しています。これにより、ユーザーは必要なリソースを独立してスケールアップまたはスケールアウトでき、コスト効率とパフォーマンスのバランスを最適化できます。また、データ共有、Time Travel、Zero-Copy Cloningなどの革新的な機能は、データ管理と分析のプロセスを大幅に簡素化します。

本仕様書では、Snowflakeの核となる概念から、具体的な使用シナリオにおけるベストプラクティスまでを網羅し、読者がSnowflakeを最大限に活用するための手助けとなることを目指します。

1. Snowflakeアーキテクチャ

Snowflakeのアーキテクチャは、従来の共有ディスクデータベースアーキテクチャと非共有データベースアーキテクチャのハイブリッドであり、データウェアハウスのすべてのコンピューティングノードからアクセス可能な永続データ用の中央データリポジトリを使用します。同時に、大規模並列処理（MPP）コンピューティングクラスタを使用してクエリを処理し、クラスタ内の各ノードがデータセット全体の一部をローカルに格納します。このアプローチにより、共有ディスクアーキテクチャのデータ管理の簡易性に加えて、シェアードナッシングアーキテクチャのパフォーマンスとスケールアウトのメリットが得られます。

Snowflakeのユニークなアーキテクチャは、以下の3つの重要なレイヤーで構成されています。

1.1. データベースストレージ

Snowflakeにデータがロードされると、Snowflakeはそのデータを内部の最適化された圧縮された列指向形式に再編成し、クラウドストレージに保存します。Snowflakeは、このデータの保存方法のすべての側面（組織、ファイルサイズ、構造、圧縮、メタデータ、統計など）を管理します。Snowflakeによって保存されたデータオブジェクトは、顧客が直接表示したりアクセスしたりすることはできません。Snowflakeを使用して実行されるSQLクエリ操作でのみアクセスできます。

1.2. クエリ処理

クエリの実行は、処理層で実行されます。Snowflakeは、「仮想ウェアハウス」を使用してクエリを処理します。各仮想ウェアハウスは、クラウドプロバイダーからSnowflakeによって割り当てられた複数のコンピューターノードで構成されるMPPコンピュータークラスターです。各仮想ウェアハウスは、コンピューティングリソースを他の仮想ウェアハウスと共有しない独立したコンピュータークラスターであり、他の仮想ウェアハウスのパフォーマンスに影響を与えません。

1.3. クラウドサービス

クラウドサービスレイヤーは、Snowflake全体のアクティビティを調整するサービスのコレクションです。これらのサービスは、ログインからクエリディスパッチまでのユーザーリクエストを処理するために、Snowflakeのさまざまなコンポーネントをすべて結び付けます。クラウドサービスレイヤーは、クラウドプロバイダーからSnowflakeによってプロビジョニングされたコンピューティングインスタンスでも実行されます。このレイヤーで管理されるサービスには、認証、インフラストラクチャ管理、メタデータ管理、クエリの解析および最適化、アクセス制御などがあります。

これらのレイヤーの分離により、Snowflakeは高い柔軟性とスケーラビリティを実現しています。例えば、ストレージとコンピューティングリソースを独立してスケールできるため、ユーザーは必要なリソースを必要な時にのみ利用し、コストを最適化できます。

2. Snowflakeの主な機能

Snowflakeは、その独自のアーキテクチャと革新的な機能により、データウェアハウスおよびデータ分析の分野で優れたパフォーマンスと柔軟性を提供します。以下に、Snowflakeの主要な機能の一部を詳述します。

2.1. セキュリティ、ガバナンス、およびデータ保護

Snowflakeは、データのセキュリティとガバナンスを最優先しており、堅牢な機能を提供しています。

- **自動データ暗号化:** Snowflakeは、保存されているすべてのデータを自動的に暗号化します。これにより、データが常に保護され、セキュリティ侵害のリスクが低減されます。
- **Snowflake Time Travel:** この機能により、ユーザーは過去の任意の時点のデータにアクセスできます。誤って削除されたデータや変更されたデータを復元したり、過去のデータの状態を分析したりすることが可能です。Time Travelの保持期間は、エディションによって異なりますが、最大90日まで設定できます。
- **Snowflake Fail-safe:** 災害復旧のために、SnowflakeはTime Travelの期間を超えて履歴データを保持します。これにより、予期せぬ障害が発生した場合でも、データが失われることなく復旧できることが保証されます。
- **列レベルおよび行レベルのセキュリティ:** Enterprise Edition以上のSnowflakeでは、特定の列や行に対してアクセス制限を設けることができます。これにより、機密性の高いデータへのアクセスを厳密に制御し、データプライバシーを強化できます。
- **オブジェクトタグ:** オブジェクトにタグを適用することで、機密データの追跡やリソース使用状況の管理が容易になります。
- **差分プライバシー:** 標的型プライバシー攻撃からデータを保護するための機能で、データの匿名性を保ちながら分析を可能にします。

2.2. データのインポートおよびエクスポート

Snowflakeは、さまざまな形式のデータを効率的にロードおよびアンロードするための柔軟な機能を提供します。

- **一括ロードおよびアンロード:** CSV、TSV、JSON、Avro、ORC、Parquet、XMLなど、多様なデータ形式に対応しています。クラウドストレージ（Amazon S3、Google Cloud Storage、Microsoft Azure）やローカルファイルからのデータロードが可能です。
- **Snowpipe:** 継続的なデータロードを可能にするサービスで、新しいデータファイルがステージに到着すると自動的にSnowflakeにロードされます。これにより、リアルタイムに近いデータ分析が実現されます。

2.3. 複製およびフェールオーバー

Snowflakeは、異なるリージョンやクラウドプラットフォーム間でのデータ複製とフェールオーバーをサポートし、ビジネス継続性と災害復旧の能力を強化します。

- **アカウント間複製:** 同じ組織内の複数のSnowflakeアカウント間でオブジェクトを複製し、データとオブジェクトの同期を維持します。これにより、データの可用性と一貫性が向上します。
- **フェールオーバー構成:** 災害発生時に、プライマリアカウントからセカンダリアカウントへのフェールオーバーを構成できます。これにより、システムのダウンタイムを最小限に抑え、ビジネスの継続性を確保します。

これらの機能は、Snowflakeが単なるデータウェアハウスではなく、包括的なデータプラットフォームとして機能することを可能にしています。ユーザーはこれらの機能を活用することで、データの管理、分析、および保護を効率的かつ安全に行うことができます。

3. Snowflakeのベストプラクティス

Snowflakeを最大限に活用し、パフォーマンス、セキュリティ、コスト効率を最適化するためには、いくつかのベストプラクティスを遵守することが重要です。以下に、主要なベストプラクティスを詳述します。

3.1. アクセス制御のベストプラクティス

Snowflakeのアクセス制御は、ロールベースのアクセス制御（RBAC）モデルに基づいており、セキュリティと管理の簡素化を実現します。効果的なアクセス制御を実装するためのベストプラクティスは以下の通りです。

- **ACCOUNTADMINロールの厳格な管理:** ACCOUNTADMINロールは、システム内で最も強力なロールであり、アカウントレベルのすべての操作を実行できます。このロールは、組織内の限られた信頼できるユーザーにのみ割り当て、日常的な操作には使用しないようにします。また、ACCOUNTADMINロールを持つユーザーには多要素認証（MFA）を義務付け、少なくとも2人のユーザーに割り当てることで、パスワード紛失時のリスクを軽減します。
- **ACCOUNTADMINロールでのオブジェクト作成の回避:** ACCOUNTADMINロールを使用してデータベースオブジェクトを作成することは避けるべきです。代わりに、組織のビジネス機能に沿ったロール階層を作成し、SYSADMINロールまたはその下位

のロールを使用してオブジェクトを作成します。これにより、最小権限の原則が適用され、セキュリティが向上します。

- **自動スクリプトでのACCOUNTADMINロールの使用回避:** 自動化スクリプトには、ACCOUNTADMIN以外のロールを使用することを推奨します。SYSADMINロールまたは階層内の下位のロールを使用して、ウェアハウスおよびデータベースオブジェクト操作を実行します。ユーザーやロールの作成・変更が必要な場合は、SECURITYADMINロールまたは十分な権限を持つ別のロールを使用します。
- **データベースオブジェクトへのアクセス権限の付与:** データベースオブジェクト（テーブル、関数、ファイル形式など）にアクセスするには、特定のオブジェクトに対する権限に加えて、コンテナデータベースおよびスキーマに対するUSAGE権限をユーザーに付与する必要があります。

3.2. データ取り込みのベストプラクティス

効率的かつ信頼性の高いデータ取り込みは、Snowflakeのパフォーマンスを最大化するために不可欠です。

- **ファイルサイズの最適化:** データをロードする際、ファイルサイズは5GBを超えないように推奨されます。ファイルサイズが大きすぎると、エラー処理が複雑になり、並列処理の利点が損なわれる可能性があります。
- **Snowpipeの活用:** 継続的なデータロードにはSnowpipeを使用します。これにより、データが到着するとほぼリアルタイムでSnowflakeにロードされ、手動での介入が不要になります。
- **エラー処理の考慮:** データロード中に発生する可能性のあるエラーを適切に処理するための戦略を立てます。Snowflakeは、エラー処理のための様々なオプションを提供しています。

3.3. ウェアハウスに関する考慮事項

仮想ウェアハウスはSnowflakeのクエリ処理の中心であり、その設定と管理はパフォーマンスとコストに直接影響します。

- **ウェアハウスサイズの選択:** クエリの複雑さとデータ量に基づいて適切なウェアハウスサイズを選択します。小さいウェアハウスはコスト効率が良いですが、大規模なクエリには不十分な場合があります。逆に、大きすぎるウェアハウスは不要なコストを発生させます。

- **自動サスペンドと自動レジューム:** ウェアハウスの自動サスペンドと自動レジューム機能を有効にすることで、使用されていないときにウェアハウスが自動的に停止し、クエリが実行されると自動的に再開されるため、コストを節約できます。
- **マルチクラスタウェアハウスの利用:** 同時実行性の高いワークロードや、異なるワークロードを分離する必要がある場合は、マルチクラスタウェアハウスを検討します。これにより、複数の仮想ウェアハウスが同時に稼働し、パフォーマンスのボトルネックを解消できます。

これらのベストプラクティスを適用することで、Snowflake環境のセキュリティ、パフォーマンス、およびコスト効率を最適化し、データ分析のニーズを効果的に満たすことができます。

4. 結論

本技術仕様書では、Snowflakeのアーキテクチャ、主要機能、およびベストプラクティスについて詳細に解説しました。Snowflakeは、その独自のクラウドネイティブアーキテクチャにより、データストレージとコンピューットの分離を実現し、比類のない柔軟性、スケーラビリティ、およびパフォーマンスを提供します。Time Travel、Zero-Copy Cloning、Snowpipeなどの革新的な機能は、データ管理と分析のプロセスを簡素化し、企業がデータからより多くの価値を引き出すことを可能にします。

また、セキュリティ、ガバナンス、およびコスト効率を最大化するためのベストプラクティスを遵守することで、Snowflake環境を最適化し、データ駆動型戦略を成功させることができます。本ドキュメントが、Snowflakeの導入、運用、および最適化に携わるすべての関係者にとって有益な情報源となることを願っています。