

PostgreSQL 解析資料

～ ソースツリーと PostgreSQL のアーキテクチャ ～

(株) NTT データ

ビジネス開発事業本部 システム方式 BU

井久保 寛明

1. はじめに

本資料では、PostgreSQL のソースツリーについて紹介する。実際にソースコードの調査を行う際に、どのあたりから見ていけばいいかという指針になることを目的としている。

ソースコードを見ていくためには、DBMS にどのような機能があるかを知っておく必要があるので、はじめに、DBMS エンジンのアーキテクチャの概要と基本的な機能を説明する。その後、ソースコード中のモジュールとの対応付けを行っていくことにする。

本資料の現状

2003年7月現在¹、PostgreSQLのソースコードの解析は、まだ、半分も終わっていないため、推測の部分が多く含まれる。そのため、誤りも少なからず含まれていることと思う。この点については、あらかじめ御了承頂きたい。

本資料は、*PostgreSQL 7.3.3* をベースに書いてあるので、他のバージョンでは、異なる場合があるので注意して頂きたい。

2. PostgreSQL のアーキテクチャ概論と機能の概要

PostgreSQL のソースツリー及びソースコードを見ていくにあたって、PostgreSQL にどのような機能があるのかを知っておく必要がある。この章では、PostgreSQL のプロセスのアーキテクチャと DBMS サーバ(backend) のアーキテクチャについて紹介する。DBMS のアーキテクチャや要素技術を押さえておくことで、ソースツリーのディレクトリ名やファイル名などから、その内容が想像できるようになる。

2.1. PostgreSQL のアーキテクチャ

まずは、PostgreSQL のプロセスに関するアーキテクチャを紹介する。PostgreSQL のクライアント / サーバ間のアーキテクチャは図 1 のようになっている。

¹ 2003年10月に改訂。9月のJPUG合宿などの情報をフィードバックしました。

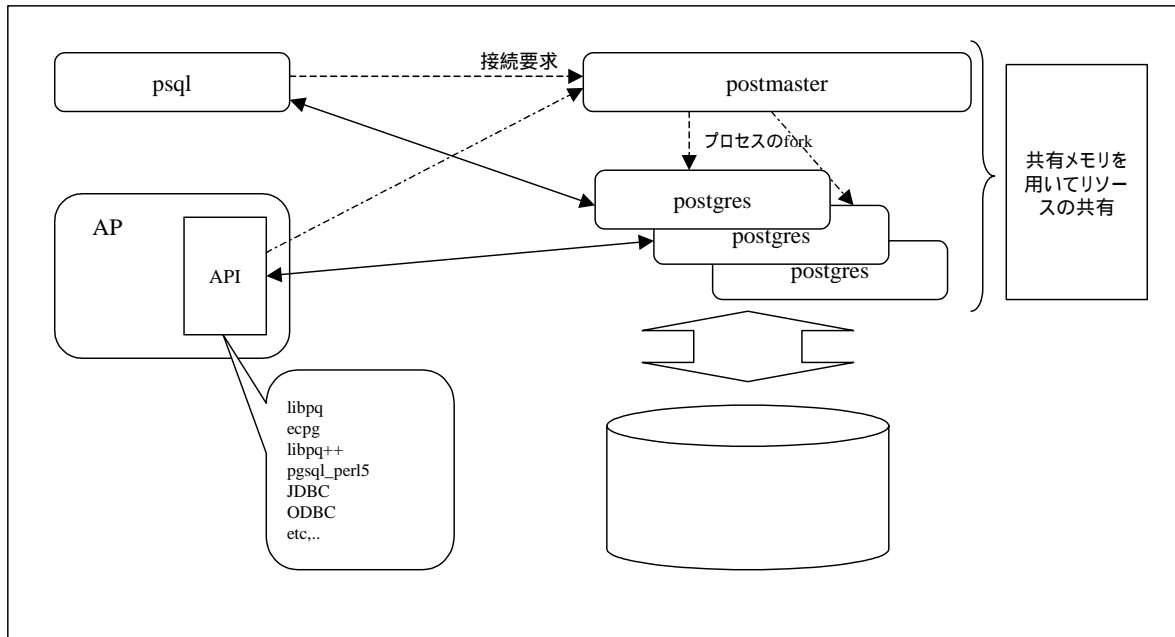


図.1 PostgreSQL のアーキテクチャ

左側 2 つがクライアントである。"psql"は、PostgreSQL 付属の SQL コマンドラインインタプリタである。DB のメンテナンスなどのときに psql から SQL 文を実行したり、シェルスクリプト上から SQL を実行したりする場合に利用できる。"AP"は、ユーザの作成したアプリケーションである。さまざまな API(Application Programming Interface)が用意されており、これらを通じて PostgreSQL のサーバへアクセスすることになる。

図の右側が PostgreSQL のサーバとして実行されるプロセスになる。postmaster とは、PostgreSQL のサーバの接続要求をまとめて受け付けるプロセスである。まずクライアントは、postmaster に接続要求を出す。すると、postmaster は 1 つの接続につき 1 つの postgres プロセスを fork する。以降、そのクライアントの処理は直接 postgres プロセスが受け持つことになる。各 postgres プロセスは、接続が切れる際に終了する。postgres プロセス同士は、共有メモリを利用して、バッファの共有やロック情報の共有を行う。

DB のデータの読み書きを行うのは、必ず postgres プロセスである。図は、クライアントサーバの例であるが、スタンドアロンで DB への処理を実行する場合、コマンドラインから postgres 自体を直接起動することができる。例えば DB の初期化スクリプト(initdb)の中で、このような使い方をしている。

2.2. DBMS エンジンのアーキテクチャ

次に、DBMS のエンジンに関する一般的なアーキテクチャを紹介する。本来なら、PostgreSQL のアーキテクチャをそのまま書きたいところだが、調査が十分に進んでいないため、ここでは一般的な DBMS のアーキテクチャを中心に紹介する。その後、3 章の中で、PostgreSQL のソースコード中のどのモジュールが、DBMS のアーキテクチャのどのあたりの処理に相当するのかを照らし合わせていく。

DBMS の機能を大きく分けると次の 3 つに分けることができる。もちろんいろいろな分け方が可能だが、本資料ではこれで説明していく。

- ストレージ管理
- クエリの実行(クエリコンパイラ、クエリエグゼキュータ)
- トランザクション管理

これらの関係を表したものが、図 2 である。

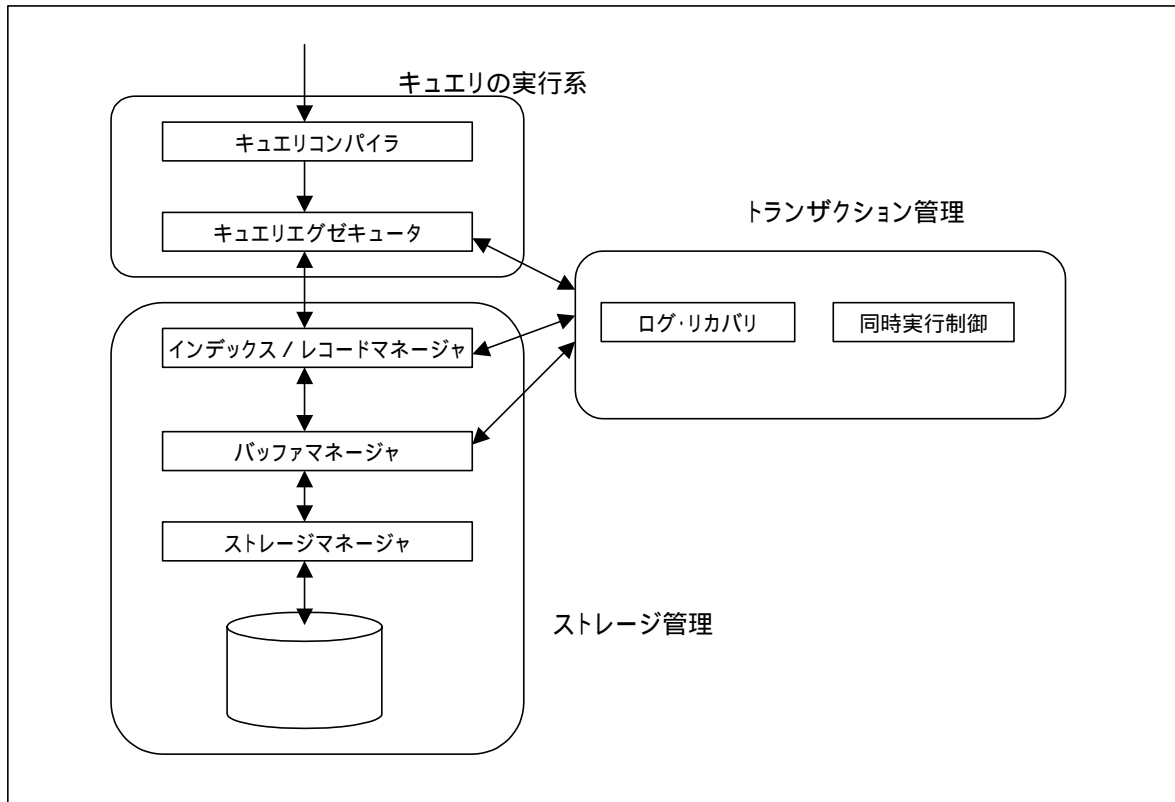


図.2 DBMS エンジンのアーキテクチャ

ストレージ管理とは、データベースに使用する OS ファイルの管理やその I/O の管理、インデックスなどの高速なファイル I/O の手段の提供、上位層へのレコード単位でのデータの供給などを担当する。基本的には OS のファイルシステムと非常によく似た技術である。

クエリの実行系は大きく 2 つの部分からなる。1 つは、クエリ処理の様々なアルゴリズムを実装してあるクエリエグゼキュータで、もう 1 つが、与えられた SQL をどのように実行するかを決定するクエリコンパイラである。OS など他の技術に例えると、クエリエグゼキュータはエミュレータのようなものであり、またクエリコンパイラはその名のとおりコンパイラである。クエリコンパイラに至っては、C コンパイラの数倍の規模になる。

トランザクションの管理では、大きく 2 つの仕事がある²。1 つは、トランザクションの原子性 (Atomicity) を管理することで、PostgreSQL ではログを使った管理を行う。ログを使った方式は、多くの商用 DBMS で用いられる主流の方法である。もう 1 つは、同時実行の管理である。PostgreSQL では、

² 正確には、ACID 特性を満たすようにする役割を担っている。Durability などは、バックアップの保持なども影響するので、ここではあえて 2 つとする。

ロックとタイムスタンプ³を使った方式を採用している。

2.2.1. ストレージ管理

最近の DBMS は、データ領域として OS のファイルを使用するものも多い。OS のファイルの中を固定長のブロック単位に分けて管理する。そして、ディスク(OS のファイル)に対しては、ブロック単位で I/O を行う。多くの場合、DBMS ではメモリより多くのデータを扱う必要があるので、ディスク上の必要なデータをメモリ中にバッファリングする。各ブロックは、メモリバッファの中ではページと呼ばれる。各ページ内には、レコード(タプル) 単位のデータが格納される。固定長のレコードの中に可変長のレコードを格納する仕組みが提供されている。

テーブルに入っているデータは次第に大きくなっていくので、固定長のブロックを集めてヒープとして管理する。

大規模なデータを扱う場合、必要なレコードをディスク上から検索するのに非常に時間がかかる。そのため、検索時間を節約するために、必要に応じてインデックスと呼ばれるデータ構造を構築する。インデックスを作成するかどうかは、ユーザに委ねられる。インデックス自体も、レコード単位のアクセスモジュールを利用して作成される。

図 3 が、これらの動作のイメージ図である。

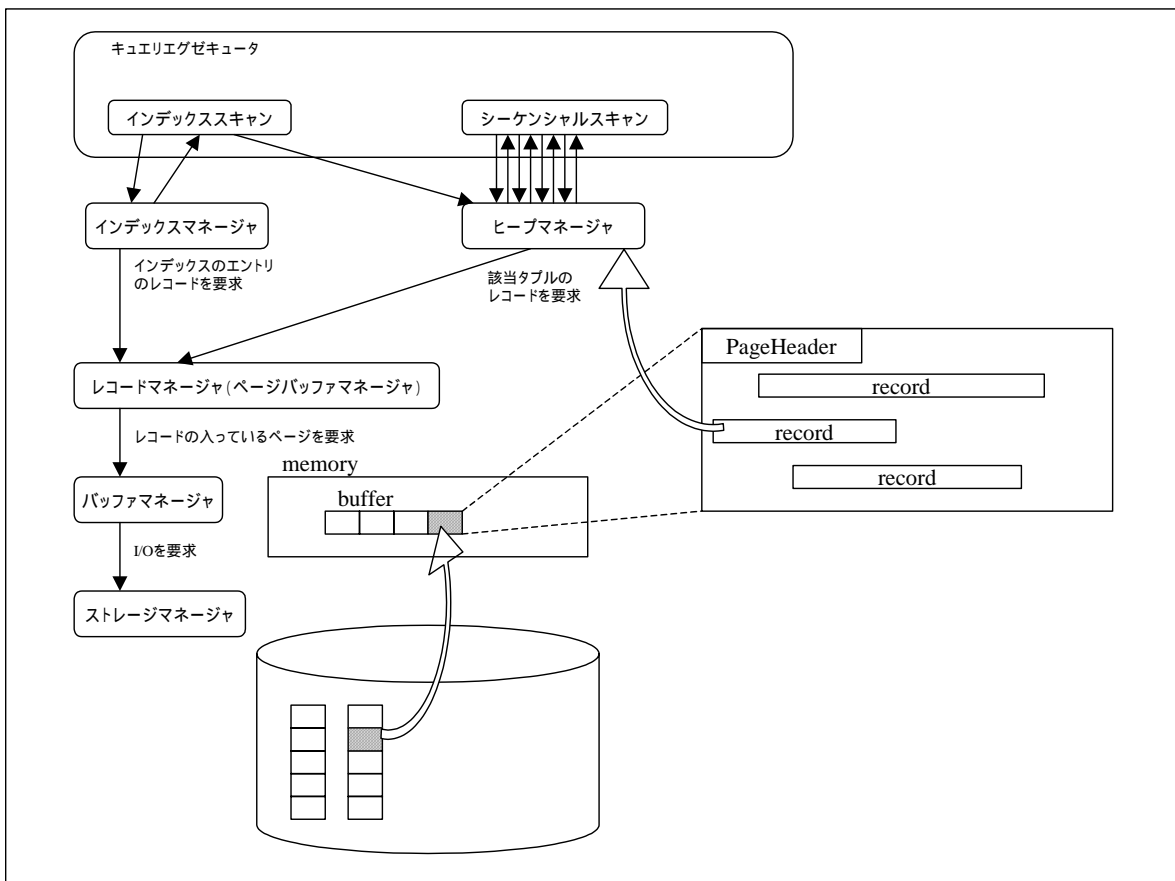


図 3. ストレージ管理

³ MVCCと呼ばれる、多版式の同時実行性制御を行う。タイムスタンプと言っても実際は、トランザクションIDを使用する。DBMSの教科書的な話では、トランザクションIDもトランザクションの順序性を管理できるので、タイムスタンプという呼び方をする。

2.2.2. クエリの実行系

クエリを実行するときの流れは、図 4 のようになる。

テキストで書かれた SQL 文は、スキャナ（字句解析器）、パーザ（構文解析器）によって**字句解析**、**構文解析**が行われて、パーズツリー（図 5 の右上のツリー）に分解される。字句解析とは、テキストを読み込んで、トークンと呼ばれる「予約語や ID、データ部分などの塊」を切り出すことである。構文解析とは、字句解析で取り出したトークンを構文規則にあてはめ、パーズツリーを生成することである。ここで、木構造にするのは、そのほうがプログラム中で扱いやすいからである。ここまでの段階で、構文のチェックなどが行われる。

続いて、**意味解析**を行う。意味解析では、テーブル名、属性名が実際に DB 内に存在するかということや、演算処理に使われる引数のデータ型のチェックなどが行われる。その他に、SQL 文で省略されているところなどを解析し、論理クエリプランへの書き換えに必要な情報を埋めておく。教科書的は意味での意味解析の処理はここまでであり、論理クエリプランの作成は後のフェーズで行われるのだが、PostgreSQL の場合は、この意味解析のフェーズでパーズツリーを論理クエリプラン（クエリツリー、図 5 の左下のツリー）に書き換える。

VIEW が存在する場合、VIEW に定義されているクエリ定義をクエリツリーに展開する。この書き換え処理を PostgreSQL では、**リライト**と呼んでいる。リライト処理では、単に VIEW の書き換えを行うだけでなく、ルールの適用も行っている。

DDL⁴、メンテナンス用のコマンドなどは、この段階で実行が可能な状態にできるので、これらのコマンド類の場合、クエリコンパイルの処理は終了して、ユーティリティ系処理の実行系に渡される。これ以降は、SQLのうちDML⁵に対して行われる。

教科書的な DBMS の話になるが、論理クエリプランは、代数学に基づいたさまざまな書き換え処理が可能な木構造である。この書き換え可能な性質を利用して、クエリの最適化を行っていく。クエリの最適化は、経験則に基づいて変換する方法や統計情報を用いて変換していく。どのような方針で最適化を行うかは、実装に依存する。

論理クエリプランの最適化が終わったら、物理クエリプランへの変換を行う。ここでは、演算のアルゴリズムの決定や中間結果の渡し方、テーブルのスキャン方法などを決定していく。物理プランでも統計情報などに基づいて最適化が必要である。これらの処理を経て、クエリプランが決定される。

PostgreSQL の場合、論理クエリプランでは、それほど最適化をかけていないように見受けられる。プランナーでは、物理クエリプラン（パスツリー、図 5 右下のツリー）を生成する。プランナーで実行可能な全ての実行パスを生成して、それぞれのパスツリーのコスト計算を行い、最適なクエリプランを選択する。

最適なクエリプランが選択できたら、そのパスツリーをクエリエグゼキュータに渡してクエリエグゼキュータを実行する。PostgreSQL の場合、パスツリーだけでなく、クエリツリーもエグゼキュータに渡している。これは、エグゼキュータの初期化に必要な情報がクエリツリーにだけ入っているからである。

⁴ DDL(Data Define Language) : スキーマやテーブルを定義するのに使用されるSQL文

⁵ DML(Data Manipulation Language) : Select, Insert, Delete, Update など、データ操作のSQL文

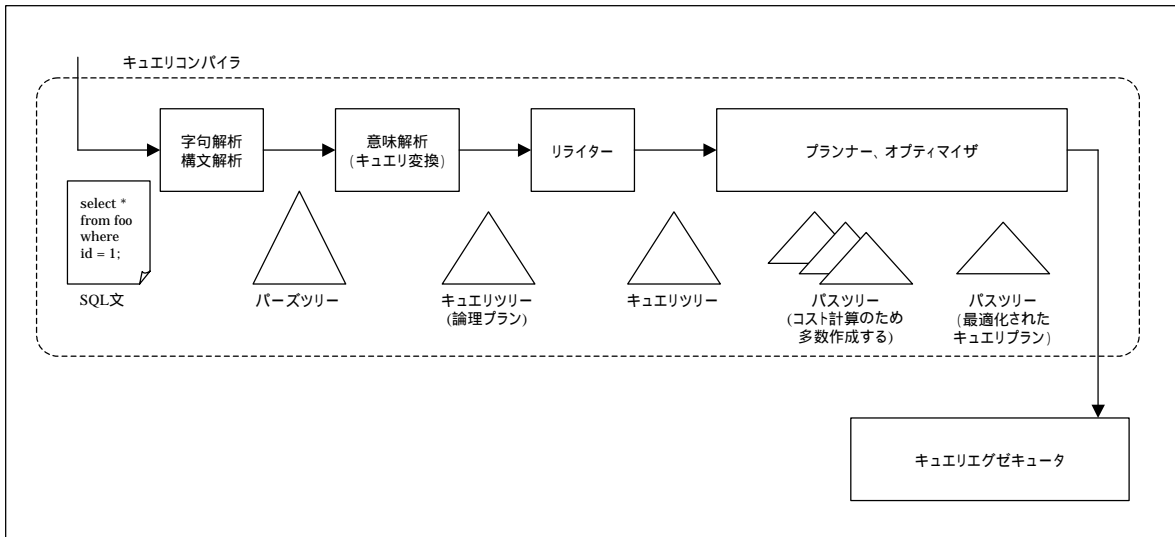


図 4. クエリの実行系

クエリエグゼキュータには、さまざまな処理を行うためのアルゴリズムが実装されている。クエリコンパイラの生成したクエリプランには、どのアルゴリズムにどの引数を渡して処理するということが細かく指示されているので、エグゼキュータは、その指示どおり動くだけである。エグゼキュータの実装では、ストレージ管理モジュールに対してレコード単位での処理要求をする。

クエリを処理していく際のデータ構造の変化は、次の図のようなイメージになる。

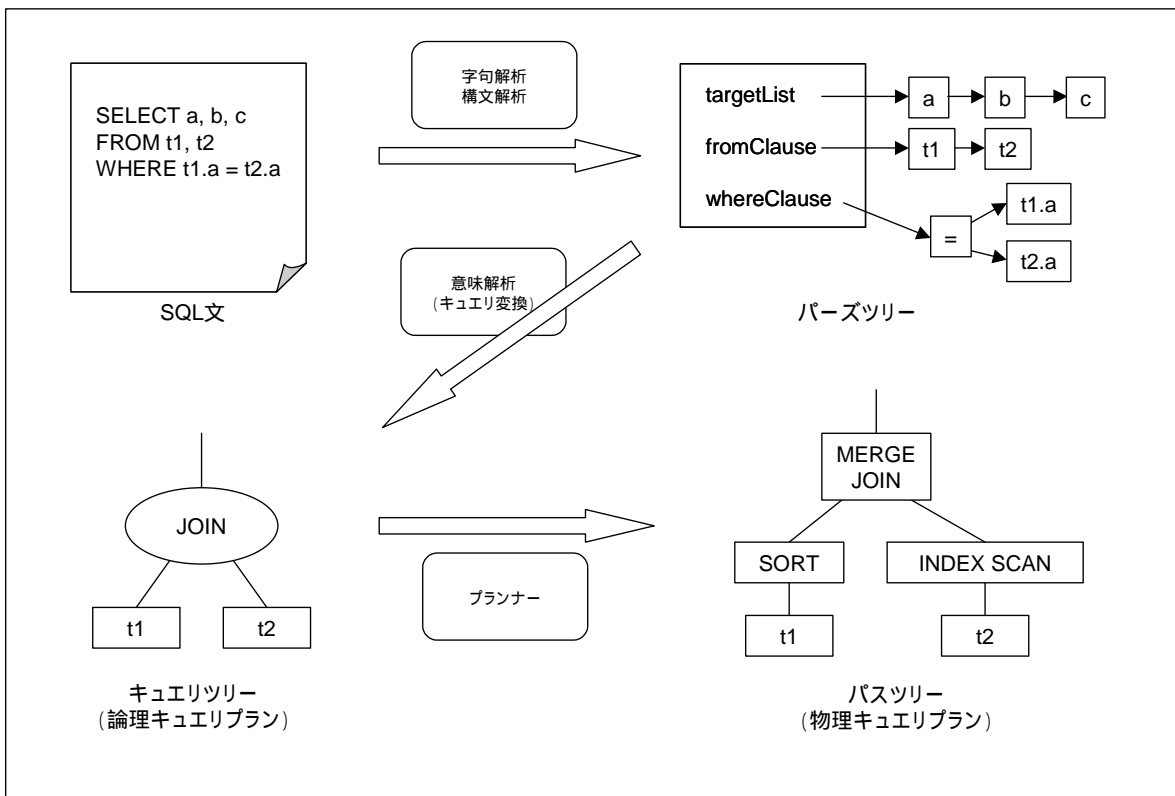


図 5 クエリのデータの変化

2.2.3. トランザクションの管理

トランザクションマネージャのログ・リカバリの機能は、データベースが壊れた場合にデータベースを復旧する機能である。ログについては、この他にもトランザクションのロールバックにも使われる。ログの内容については、データ自体の変更やデータベースの構造の変更などが書き込まれている。ログに使われるログデータバッファは、通常のデータ用のバッファとは別に管理されている。

リカバリ処理とは、ログを読みながら、個々のログエントリの処理を再実行するものである。

有名な WAL(Write Ahead Logging) とは、ログの書き込み方式で、データバッファをディスクに書き込む前に、必ずその部分のログをディスクに書く方式である。

ちなみに、PostgreSQLのログ・リカバリのモジュールは、追記型の特性(データに対するUNDOログ⁶がない)により非常にシンプルに書かれている。

PostgreSQLの同時実行制御は、ロックとタイムスタンプを併用している。タイムスタンプと言っても、実際はトランザクション ID を使用する。DBMSの教科書的な話では、トランザクション ID もトランザクションの順序性を管理できるので、タイムスタンプという呼び方をする。PostgreSQLでは、タイムスタンプと追記型のストレージの特性を利用して MVCC(多版型同時実行制御)を実現している。これは、あるレコードが更新中の場合、別のトランザクションは更新前の値を読み込むことで同時実行性を向上させる機能である。

ロック管理のモジュールはロックテーブルを管理している。ロックテーブルを用いるロックの他に、スピンロックも実装されている。

⁶ データを更新前の状態に戻すためのログのこと。データを更新後の状態にするためのログをREDO ログと呼ぶ。

3. PostgreSQL のソースツリー

PostgreSQL 7.3.3 のソースツリー直下には、次のようなファイルが含まれる。

```
poseidon(2)% cd postgresql-7.3.3
poseidon(3)% ls -F
COPYRIGHT      INSTALL      acllocal.m4   configure.in   register.txt
GNUmakefile.in Makefile     config/       contrib/      src/
HISTORY        README      configure*    doc/
poseidon(4)%
```

コンパイルする際は、ここから、`configure7` を実行し、`make` することになる。詳しくは `INSTALL` ファイルを参照して頂きたい。

ソースコードについては、`src` の下にある。本資料では、この下の特に `backend` に絞って説明していく。

3.1. src 配下

以下のファイル、およびディレクトリが `src` 配下に含まれるものである。

DEVELOPERS		FAQ が <code>pgsql/doc/FAQ_DEV</code> にあり、開発支援ツールが <code>pgsql/src/tools</code> にあるというメモ
Makefile		各サブディレクトリに対して、 <code>make</code> の実行を引き継がせる定義が書いてあるだけ。
Makefile.global		<code>make</code> の <code>suffix</code> などのルールを書いたファイル。Makefile などから読み込まれる。 <code>configure</code> 実行時に <code>Makefile.global.in</code> から作成される
Makefile.global.in		Makefile.global の元ファイル。
Makefile.port		プラットフォームごとの定義の書かれた <code>makefiles</code> 以下のファイルへのシンボリックリンク。Makefile.global から <code>include</code> される。 <code>configure</code> 実行時に作成される。
Makefile.shlib		Shared ライブラリ作成のルールが書かれた Makefile。Shared ライブラリが使えるプラットフォームの Makefile で <code>include</code> される。
nls-global.mk		NLS のルールが定義されているファイル。Makefile.global から <code>include</code> される。
win32.mak		Windows 用の <code>makefile</code> 。ココには、ODBC などでクライアントとしてアクセスするのに必要な部分しか含まれていない。
backend/		PostgreSQL のバックエンドの本体が入っているディレクトリ

⁷ ソースコードをコンパイルする際に、一番はじめに実行してコンパイル環境に合わせた Makefile などを作成するコマンド。

bin/		インストール後の bin ディレクトリ以下に納められるツール類のスクリプトやソースを格納したディレクトリ。
corba/		CORBA I/F
data/		charset の変換テーブルのサンプル (日本語では使えない?)
include/		C でクライアントを書く際に必要な include ファイルを集めたディレクトリ
interfaces/		各種 API の実装。ecpg, jdbc, libpgtcl, libpq, libpq++, odbc, perl5, python などが実装されている。
makefiles/		Makefile 中の各プラットフォームの固有の定義。Makefile.port としてシンボリックリンクを張って、Makefile.global に読み込まれる。
pl/		Stored Procedure 用の言語の実装。
template/		各種プラットフォーム用のコンパイルの引数のテンプレート
test/		コンパイル後のテストのためのコード類
tools/		PostgreSQL 開発者向けのツール
tutorial/		チュートリアルとしてのサンプルコード
utils/		ユーティリティ。現在、DLL 初期化ユーティリティのみ。

src 配下の主なディレクトリは、backend, include, interfaces である。

backend が PostgreSQL のサーバにあたる部分のコードである。これについては、後述する。

include には、ソースコードのうちアプリケーションプログラムでもインクルードする必要があるヘッダファイルが入っている。当然、backend, interfaces 配下のソースでも使用しているので、ソースコードを解析する際にはここも見る必要がある。

interfaces には、図 1 にも書いたような各種 API のソースが入っている。

ecpg, jdbc, libpgtcl, libpq, libpq++, odbc, perl5, python などである。

libpq については、backend/libpq というものもあるが、ライブラリとしての libpq は interfaces 配下にある方である。backend 配下の libpq は、libpq プロトコル⁸の backend 側の実装になる。

3.2. backend

以下の表が、backend ディレクトリ配下に含まれるモジュールである。

Makefile		
1 access/		
	1.1 common/	タプル操作、タプルのヘッダ操作の補助関数などが書いてある。
	1.2 gist/	GiST index(多次元インデックス)の実装。
	1.3 hash/	ハッシュインデックスの実装。
	1.4 heap/	heap の実装。catalog 配下の heap.c から呼ばれる。こちらが実ファイルアクセスに近い側。
	1.5 index/	index のインタフェース部分の実装。
	1.6 nbtree/	Btree インデックスの実装。

⁸ libpq ライブラリを用いて通信する方法を PostgreSQL では、特に libpq プロトコルと呼んでいる。

	1.7 rtree/	2次元インデックス ⁹ R-treeの実装。
	1.8 transam/	トランザクションマネージャ(ログ管理部分)
2 bootstrap/		初期データベース構築時のモジュール
3 catalog/		カタログの実装
4 commands/		エグゼキュータのうち、コマンド(DML 以外の処理)部分の実装
5 executor/		エグゼキュータのうち、クエリエグゼキュータの実装
6 lib/		内部で使用するユーティリティ的なライブラリ。
7 libpq/		サーバ側の libpq プロトコルの実装。主にユーザ認証など
8 main/		postmaster, postgres などのプログラムの main() 関数の実装
9 nodes/		SQL コンパイラで使用するノードオブジェクトの実装
10 optimizer/		SQL コンパイラのうち、オプティマイザの実装
	10.1 geqo/	Genetic Query Optimization (遺伝的クエリ最適化)部分の実装
	10.2 path/	JOIN するときに、可能な全ての順序のパターンを生成する
	10.3 plan/	クエリプランの生成を行う
	10.4 prep/	特殊なケースに対するプリプロセスフェーズの実装
	10.5 util/	オプティマイザ用の補助関数などの実装
11 parser/		クエリコンパイラのうちスキャナ、パーザ、意味解析部分の実装
12 po/		メッセージカタログのようなもの。全体としてあまり使われていない?
13 port/		各プラットフォームの移植用コード
14 postmaster/		Postmaster の実装
15 regex/		正規表現関連の実装
16 rewrite/		SQL コンパイラのうち、リライターの実装
17 storage/		ストレージ全般を管理するモジュールを集めたディレクトリ
	17.1 buffer/	バッファ管理。ここで管理するのは、データ部分及びインデックスに使用するバッファであり、ログ用のバッファは、別に用意される。
	17.2 file/	仮想ファイルデスクリプタの実装。
	17.3 freespace/	ページ内 (テーブル内) の空スペースを管理する。
	17.4 ipc/	共有メモリの管理
	17.5 large_object/	blob など lo の実装。
	17.6 lmgr/	ロックマネージャ。
	17.7 page/	ページ中のデータ管理モジュール。レコード単位のデータの取り出しを実装。
	17.8 smgr/	実ストレージへの I/O を実装したモジュール。ディスクへのシステムコールなどを管理。
18 tcop/		PostgresMain などの実装が入っているところ。
19 utils/		
	19.1 adt/	データ型に関する実装。

⁹ 地図データのようなものを扱わない限り、普通のトランザクション処理で使用されることはまずない。

19.2 cache/	カタログデータ、リレーション、ファンクションのキャッシュを担当するモジュール。
19.3 error/	エラーメッセージ出力に関する実装
19.4 fmgr/	ファンクションマネージャ: ファンクションの実装。
19.5 hash/	ハッシュ関数の実装
19.6 init/	postgres または postmaster 起動時の各種初期化モジュール
19.7 mb/	マルチバイト処理に関するモジュール
19.8 misc/	実行時のパラメータを扱う GUC ¹⁰ やその他汎用モジュールがある場所。
19.9 mmgr/	メモリマネージャ。階層的にメモリを管理できるモジュールを提供。
19.10 sort/	ソートの実装。
19.11 time/	MVCC の実装 (タイムスタンプ型の同時実行制御)

図 2 のクエリコンパイラにあたる部分のモジュールが、9 nodes, 10 optimizer, 11 parser, 16 rewriter である。11 parserは、SQL(その他のPostgreSQLのコマンドも含む) の字句解析¹¹を行って構文解析¹²を行うモジュールと、意味解析¹³を行うモジュールからなる。

図 2 のクエリエグゼキュータにあたるのが、4 commands と 5 executor である。

図 2 のトランザクションマネージャのうち、ログ・リカバリにあたるのが、1.8 transam 配下のモジュールである。また、ログ・リカバリのためのコードの一部は、ロギングの必要な各モジュールにも埋め込まれている。図 2 のトランザクションマネージャのうち、同時実行制御の部分が、17.6 lmgr と 19.11 time である。

図 2 のストレージ管理のインデックス/レコードマネージャにあたるのが、1 access 以下の transam 以外のモジュールと 17.7 page である。バッファマネージャは 17.1 buffer あたりで、ストレージマネージャは、17.2 file や 17.8 smgr あたりになる。

その他、図 2 の構造に出てこないもので、主要なモジュールには、次のようなものがある。

2 bootstrap:

データベースを初期化する際に使用するモジュールが含まれている。主にシステムカタログのデータを作成するために使用される。データベースの初期化時に読み込むファイルは、C のヘッダファイルからソースコードコンパイル時に自動生成される。ヘッダファイルから自動生成したファイルを読み込むことによって、ソースコード中に埋め込まれている OID と、データとしてシステムカタログに取り込まれる OID を一致させるための機構である。ファイルを読み込む機能のため、独自構文を利用している。しかし、SQL に比べてはるかに簡単なものなので、実装はきわめて小さい。

¹⁰ Grand Unified Configuration Scheme の略。コンフィグで設定できるパラメータを扱う機構。

¹¹ SQL文からトークン(予約語やid、データ部分などの塊) を切り出す

¹² 構文規則に従って、パース木を生成する

¹³ テーブル名や属性名が実際にDBに定義されているかをチェックしたりする

3 catalog:

システムカタログを実装したモジュール。システムのメタ情報を保持しているのだが、PostgreSQL の場合、この部分を使って、データ型や関数、その他いろいろなものを定義できるようにしてある。

8 main:

PostgreSQL のサーバプログラムのエントリポイントである。PostgreSQL をコンパイルすると、postgres, postmaster というの2つのプログラムがあるが、これらはファイル名が異なるだけで同一のプログラムである。main には、これらのプログラムが最初に呼び出す main() 関数が含まれている。

14 postmaster:

main から postmaster として呼び出されたときの実装。クライアントからの接続を待ち、接続があった場合に子プロセスを fork する。

18 tcop:

main から postgres で呼び出されたとき、および postmaster から fork されたあと呼び出される PostmasterMain() の実装が入っている。

6 lib:

ライブラリとして使用するようなユーティリティ的のものが入っているディレクトリ。

7 libpq:

サーバ側の libpq プロトコルの実装。主にユーザ認証などである。

19 utils:

汎用モジュールのようにまとめてあるが、データ型に関する実装である 19.1 adt、リレーションやカタログデータなどのキャッシュを行う 19.2 cache、メモリ管理を行う 19.9 mmgr などの非常に重要なモジュールも含まれている。

4. ソースコードの流れ

ここでは、ソースコードを読み始める糸口として、Cプログラムの先頭である main()からの流れを紹介する。

PostgreSQL のサーバを上げるためのコマンド postmaster および スタンドアロンで DBMS を実行するためのコマンド postgres は、main/main.c の main()関数から始まる。

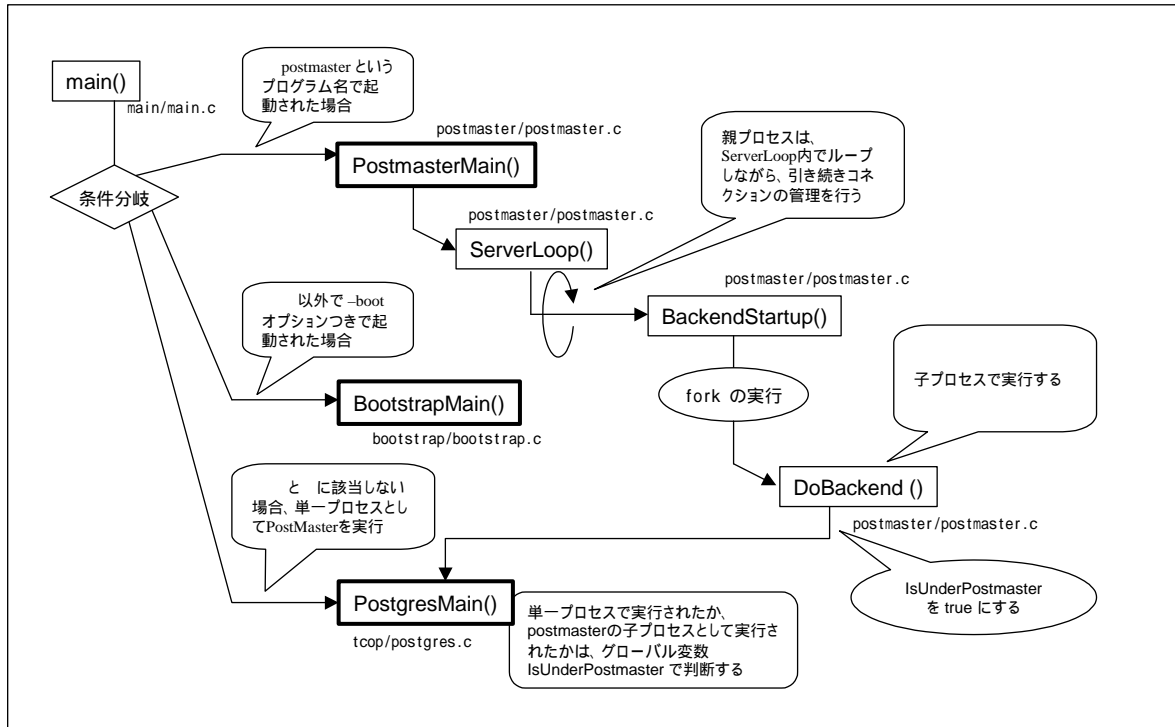


図 6. main() 関数からの流れ

main()関数の中で、postmaster という名前と呼ばれていたら、postmaster として起動する（ のルート）、-boot オプション付で呼ばれると、BootstrapMain という関数を呼び出し初期化専用プロセスとして動作する（ のルート）、その他の場合は、スタンドアロンの PostgreSQL のデータベースアクセスプログラムとして動作する（ のルート）。

コマンドラインから postgres として呼び出されたときの動作は、同時実行制御を省略していることとプロセスの初期化を行うこと以外は、postmaster から fork されている postgres と同じである。コード上はまったく同じ物で、コード内で postmaster から fork されたかどうかで、処理の分岐が若干入っている。

PostgresMain() 以降の処理は、図 7 のようになる。図 4 のクエリ処理の流れが左側に示してあり、それぞれのエントリーポイントが分かると思う。

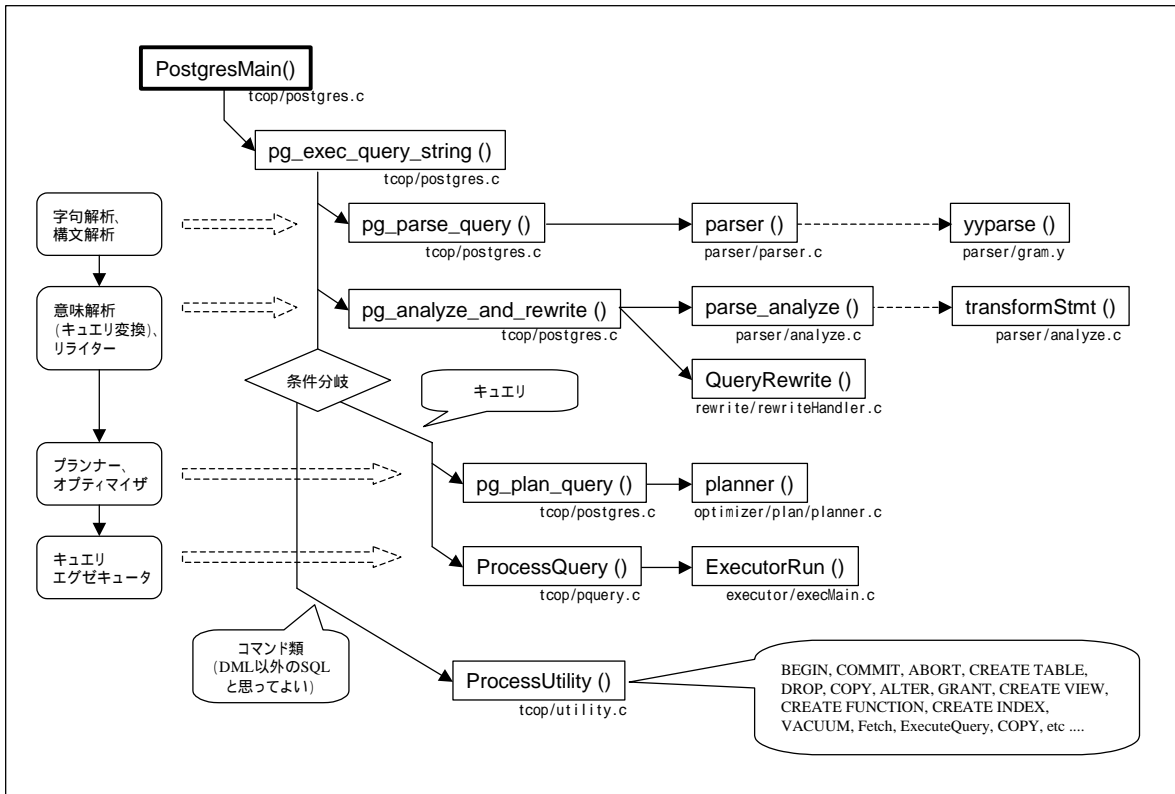


図 7. PostgresMain() 関数からの主な処理の流れ

以上が、これまでに紹介した DBMS のアーキテクチャのエントリーポイントとなる。これらの関数から PostgreSQL の backend の処理の流れを追うことができると思う。

< EOF >