



# libev+OpenSSL =EvServer (PostgreSQL編)

Future Versatile Group/MyDNS.JP  
T.Kabu

2020.11.02 #pgunconf



# ▶ ふと思いついて…

- **Postgres用マルチメインDB**の構築方法を考察  
→ PgBouncerの内部構造を解析してみた
- クライアントからの接続とデータを受け付けて  
PostgreSQLに対してリレー = PgなProxyか
- 接続 / 切断を繰り返さない = 負荷が軽いらしい

…けっこうスパゲッティで挫折しかける



# ▶ やりたいことは...



Clientからのクエリにより、取得(SELECTとか)なら最も速いDBから、  
その他変更系なら全ての対象DBにトランザクションでクエリを発行して、  
全結果によりコミットかロールバックが決める、みたいなのでどうよ？

(もちろん処理の内容によりこの辺は細かく制御しないよね)



# ▶ すでにないの？



↑ こいつと似たようなことをしているのが

- PgPool-II 多機能
- PgBouncer シンプル&高速



# ▶ お金払えば...



まさにPostgreSQLでマルチメインできるのが

・ PostgreSQL BDR

(1ライセンス\$10,000~とか)



# ▶ PgBouncer

- とりあえずソースをGitHubで眺める
- C言語で記述されてる(なんとか読める!)
- **libevent**を使用(イベントドリブンなのねー)
- 記述の統一感が!?(切った貼ったスパゲッチ!?)
- ともあれLinux上でVSCでデバッグ環境構築
- ブレークさせながらソースにコメント付けて読解



# ▶ libevent

- イベント処理を汎用的に記述させるライブラリ
- selectやepollとかをラップして移植性を上げる
- その後libuvとかlibevとかより良いものが登場
- PgBouncerで、libeventの動作概要を把握  
→ 実際ソース起こすなら改善版の **libev** だよな

…でも実際libevでSSL/TLSまでやってる例が見当たらず、libeventと同様にできるか不安



# ▶ やってから悩もう!

- PgBouncerの動作をだいたい理解したけど…
- ナウい**イベントドリブン**はやったことないから不安
- とはいえ通信の世界では状態&イベントは常識  
(イベントマトリックス、関数の多次元配列化)

とりますケルトン作ってlibevを把握しよう!!





# ▶ libevの基本

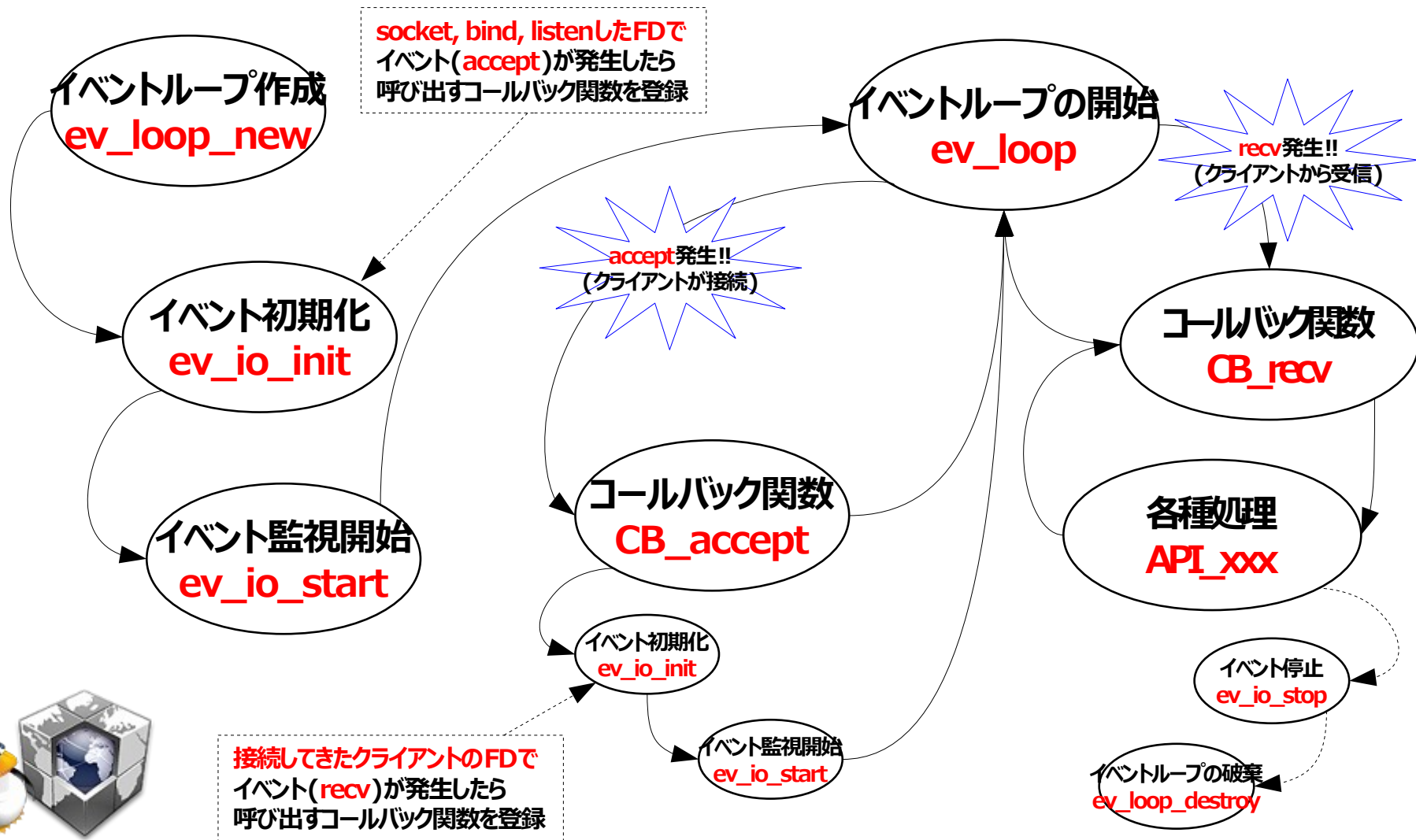
- 扱えるイベント(=割込)の種類が豊富
- イベントの処理(初期化やコールバック関数の登録)が簡単
  - 1) イベントループ作成 ⇒ `ev_loop_new`
  - 2) イベント初期化 ⇒ `ev_io_init`
  - 3) イベント監視開始 ⇒ `ev_io_start`
  - 4) イベントループの開始 ⇒ `ev_loop`
  - 5) イベントループの破棄 ⇒ `ev_loop_destroy`

割込 = イベントが発生した時に呼び出す

コールバック関数を `ev_io_init` で登録するだけ!



# ▶ libevのイメージ



# ▶ OpenSSLの基本

CentOS7はOpenSSLが1.0.x、CentOS8は1.1.x…  
(ただでさえ、みんな泣いている話しか聞かないのに)

## 1) OpenSSLの各種初期化处理

→OpenSSL 1.1.0以降、明示的な初期化は必要ないヨ

## 2) SSL設定情報を作成&設定

→許可プロトコルとかサーバー証明書とかは1.1.xでも必要

…で、普通に、**socket**→**bind**→**listen**までは同じだが、  
このあとの**accept**での接続受付からの**recv**で受信という  
のが非SSL通信とちよつと違う。



# ▶ OpenSSLの基本

まずクライアントから接続が来た(**accept**)なら…

3) **SSL\_new**でSSL接続情報を作成

4) ファイルディスクリプタを、**SSL\_set\_fd**でSSL接続情報に紐づけ(**bind**みたいなもん)

で、クライアントからデータが来た(**recv**)なら

5) **SSL\_accept()**でSSLセッションを確立させる

6) SSLセッションが確立したら**SSL\_read**でデータ受信

7) SSLセッションで必要に応じて**SSL\_write**でデータ送信

※復号化/暗号化処理は勝手にやってくれるのでラクチン!



# ▶ OpenSSLの基本

で、クライアントと切断の際には

8) **SSL\_get\_fd**でSSLセッションで使用しているFDを取得

9) **SSL\_free**でSSL接続情報を破棄

10) そしたら普通に**close**すればOK

こんな感じで、非SSLなセッションを確立したら、その上に改めてSSLなセッションを確立する感じで、大したことはない。

問題は、

**SSL\_accept**や**SSL\_read**で  
**イベント発生**するの？



# ▶ libev & OpenSSL

・結果から言うと、難しく考えなくてもうまくいった  
・「何か」データが来るとイベントが発生するので、**recv**処理でその接続がSSL/TLS接続かどうかで受信処理を分ければOK

1) 非暗号化通信なら

→そのまま受信(**read**)

2) SSL/TLS接続モードなら

→2-1) SSL/TLSハンドシェイクしてSSL/TLS接続を確立

→2-2) 接続確立後なら**SSL\_read**で復号データを受信

※ STARTTLSのような、非暗号通信から暗号通信に途中から切り替わる場合も、似たり寄ったり





# ▶ スケルトン完成…

- GitHUBで『**EvServer**』として公開中

<https://github.com/disco-v8/EvServer>

- 例によって日本語による、ほぼ一行一コメント
- 勉強するにはちょうどいい…のかな？感想歓迎！



# ▶ 例えばHTTPD

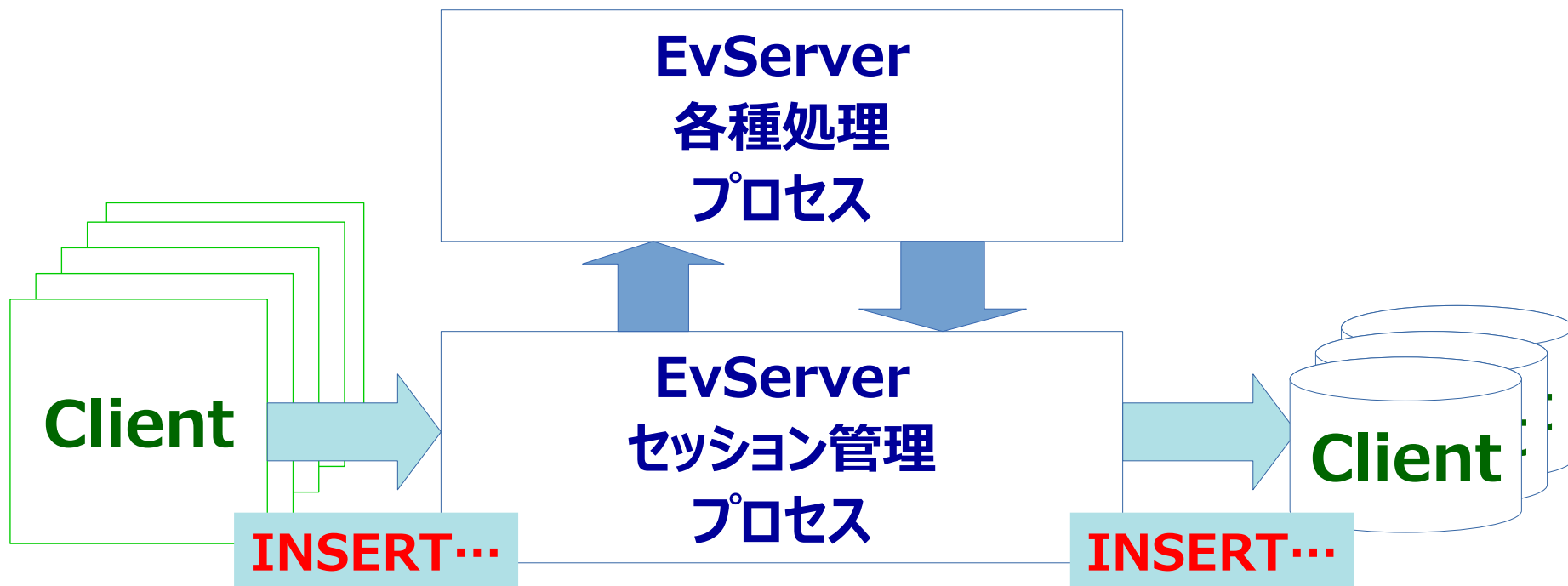
- nginxやApacheみたいな**WEBサーバー**を作成
- てけとーWEBサーバーですが、**FastCGI**も実装!
- **GET**だけでなく**POST** (multi-part)にも対応!
- Apacheよりは、**1.2~1.3倍速い**です (ab測定)
- まさか二十年以上たってからまたFastCGIを触るとは!!

というか、FastCGIプロトコル概念が古すぎだろ…ブツブツ…64kBだと!?





# ▶そして発展形(本題)



WEBサーバーは一段落したので、本来のPostgreSQL方面を。  
面倒な処理は、FastCGI同様に別プロセスとソケット通信で接続し、  
そちらに任せた方が、C言語で全部やるよりぜったいに楽かなー(笑)



# ▶ PostgreSQL

- ・『フリーなオープンソースのリレーショナルデータベース管理システム』 (<https://lets.postgresql.jp/> より)
- ・CLIからのアクセスツールとして **psql** がある  
psql hoge\_db -U hogehoge -p 5432 -h localhost
- ・psqlとPostgreSQLのやり取り(プロトコル)を把握  
→第52章 フロントエンド/バックエンドプロトコル

<https://www.postgresql.jp/document/12/html/protocol.html>



# ▶ 接続手順

- ・クライアントからの開始メッセージは二種類
  - 1) **SSLRequest** (SSL接続を要求)
  - 2) **StartupMessage** (実際の接続要求)
- ・最新のpsqlは標準でSSLRequestから
  - PostgreSQLは" S " (=OK)か" N " (=NG)を返す
- ・SならSSLハンドシェイク開始 (平文から暗号化通信へ)
- ・Nなら (もしくは暗号化通信確立後) に、2) がくる



# ▶ 接続確立後は…

- PostgreSQLから **AuthenticationOk** で接続確立
- 接続情報が **ParameterStatus** で色々送られて来る
- 最終的に **ReadyForQuery** が来たらクエリを送れる
- そしたら psql からクエリを送れるので **Query** で送ると…
- PostgreSQLからは、**RowDescription**、**DataRow** (複数の場合あり)、**CommandComplete** のように応答があり…
- 再び **ReadyForQuery** が来たらクエリを送れる



# ▶ フォーマットが...

・やり取りは簡単だけど、メッセージのフォーマットが...

1) 開始メッセージ系は最初の四バイトがメッセージ長

**XX XX XX XX** XX ... ←int32に変換する必要あり

※この長さで **SSLRequest** かどうか判定

2) それ以外のメッセージは

**XX** **XX XX XX XX** **XX XX XX XX XX** ...

..... メッセージ本体(メッセージ長-4バイト)

..... メッセージ長(4バイト、自身を含む)

..... メッセージタイプ(1バイト)



# ▶ とりまさておき...

- ・メッセージタイプはいろいろあってPostgreSQLの歴史を感じる(笑)
- ・前述のpsqlとPostgreSQLのやり取りを実際に把握しないと、本題に取り掛かれないよなあ

まずはEvServerをカスタマイズして  
やり取りを可視化するPgAnalyzerを作ってみるか!!  
まさにこれというのが↑見当たらなかったんで



# ▶ PgAnalyzer

- psqlとPostgreSQLの間に入り、メッセージを解析しつつバケツリレーするだけ

実際に動いているところを  
見ていただいた方が速いですね



# ▶ OpenSSLに泣く

- PostgreSQLではナウでヤングなsha256認証まで  
やったら、自分もOpenSSLで泣きました

実際に動いているところを  
見ていただいた方が速いですね  
(ノド` )泣泣...





# ▶マルチメインDB目指して

・簡易問い合わせ (**Query**) の中に以下の命令がある場合、

- ・INSERT
- ・UPDATE
- ・DELETE

データベースの書き換えが発生する、という事になる。

・上記以外の

- ・SELECT

については、自分から一番近い(速い)データベースに対してのみクエリを投げればいいのか？



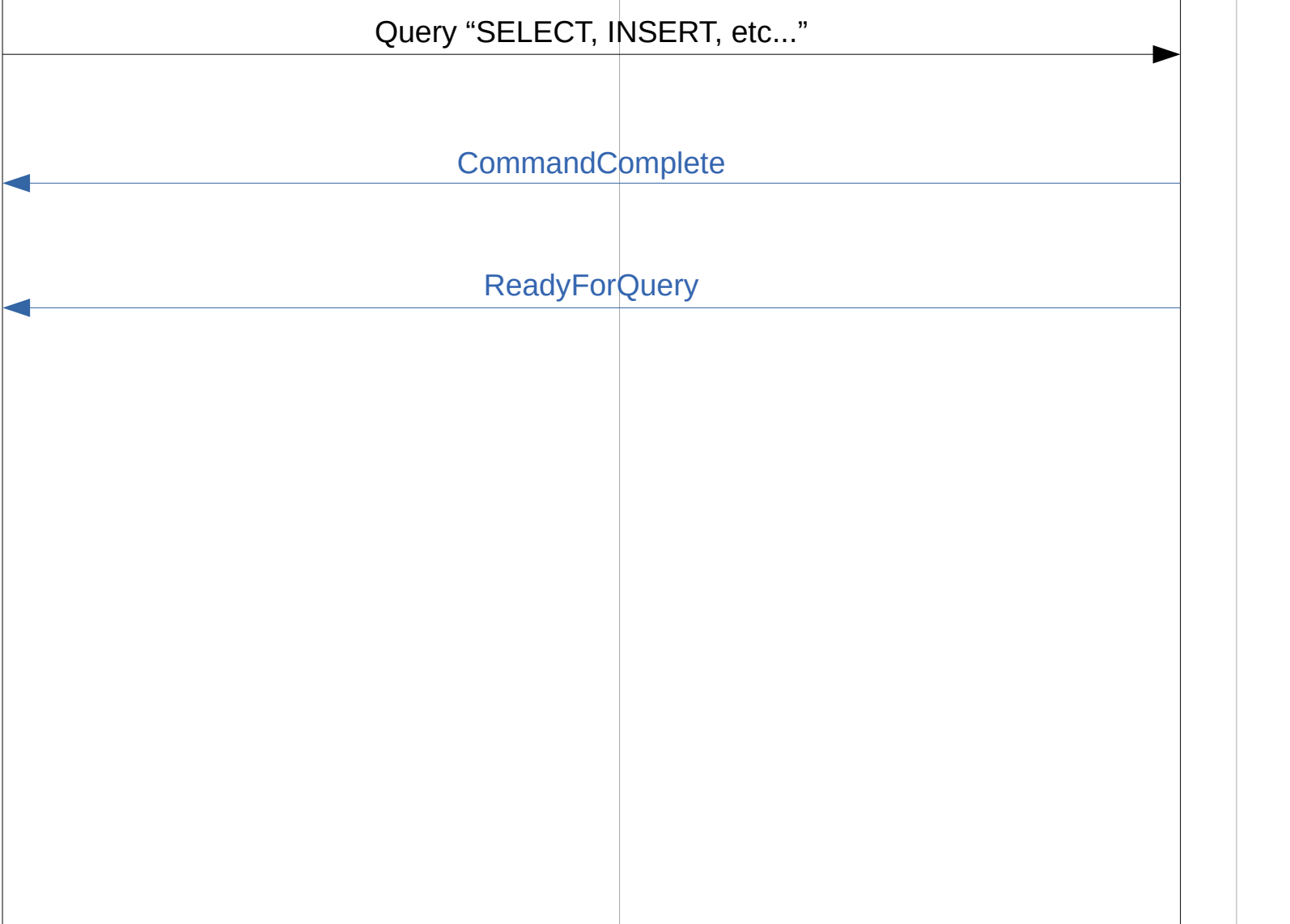
# PostgreSQL Multi DB Sequence Chart

PostgreSQL

Client

PgMulti

A B C



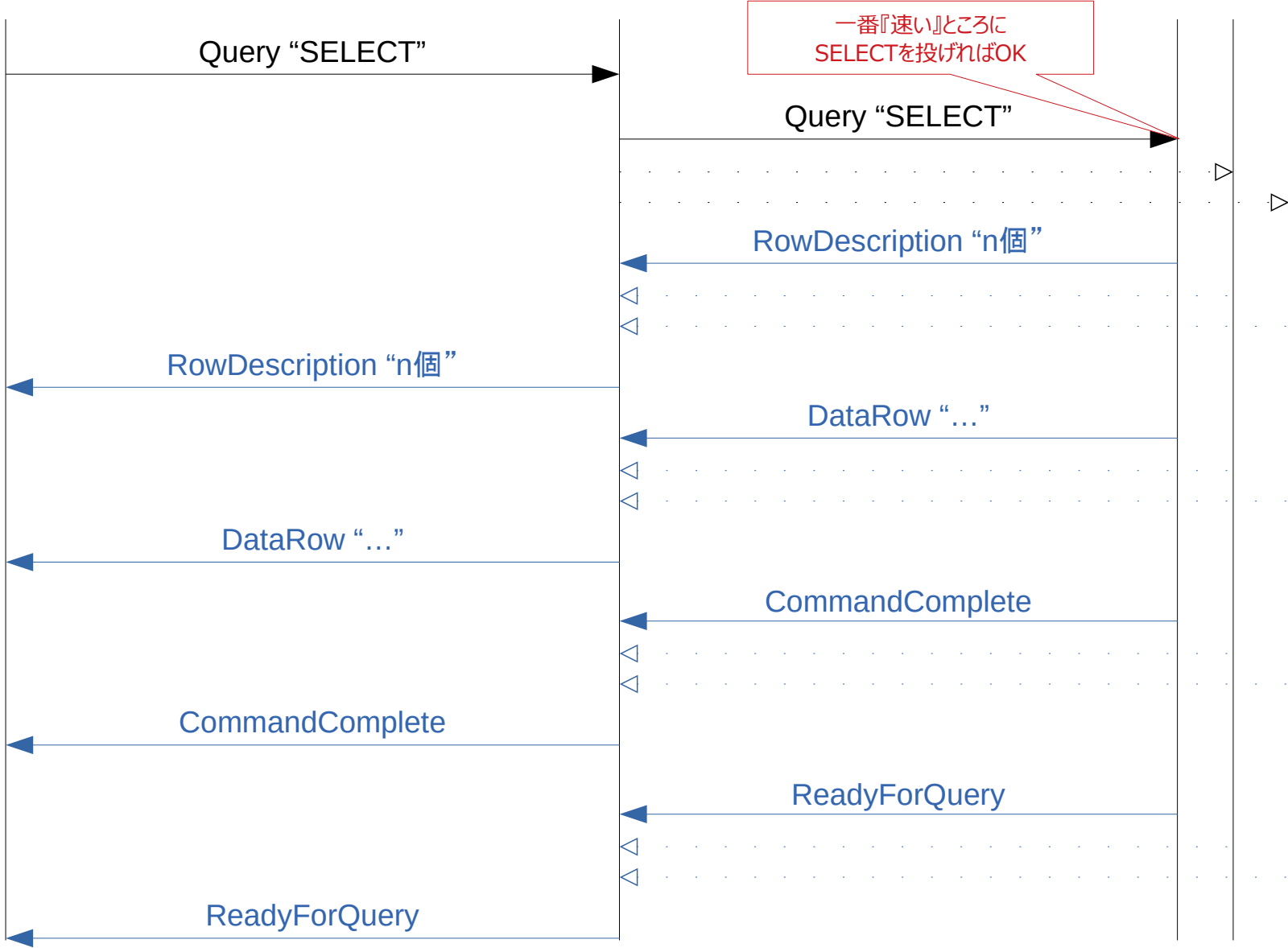
# PostgreSQL Multi DB Sequence Chart

PostgreSQL

Client

PgMulti

A B C



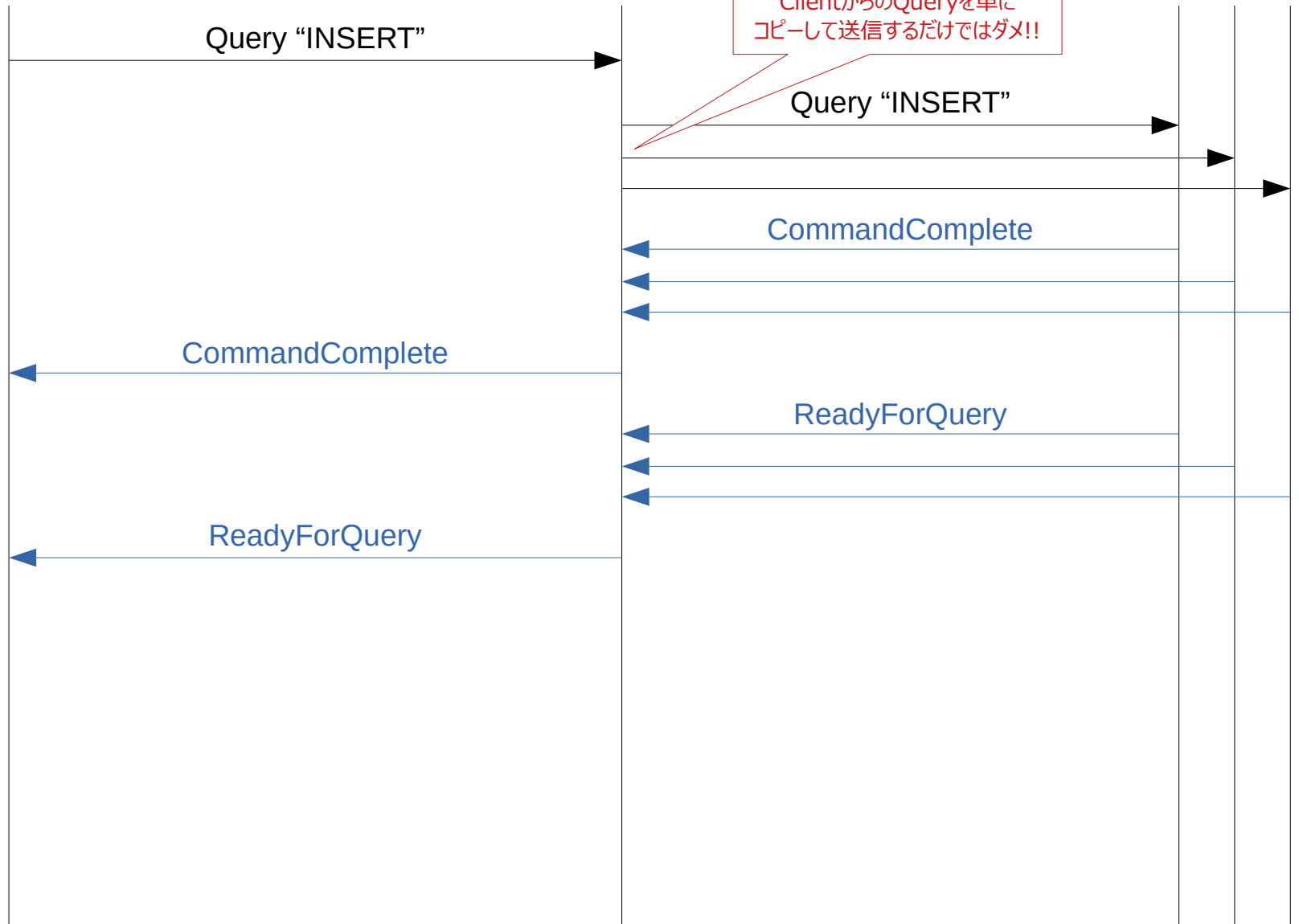
# PostgreSQL Multi DB Sequence Chart

Client

PgMulti

PostgreSQL

A B C





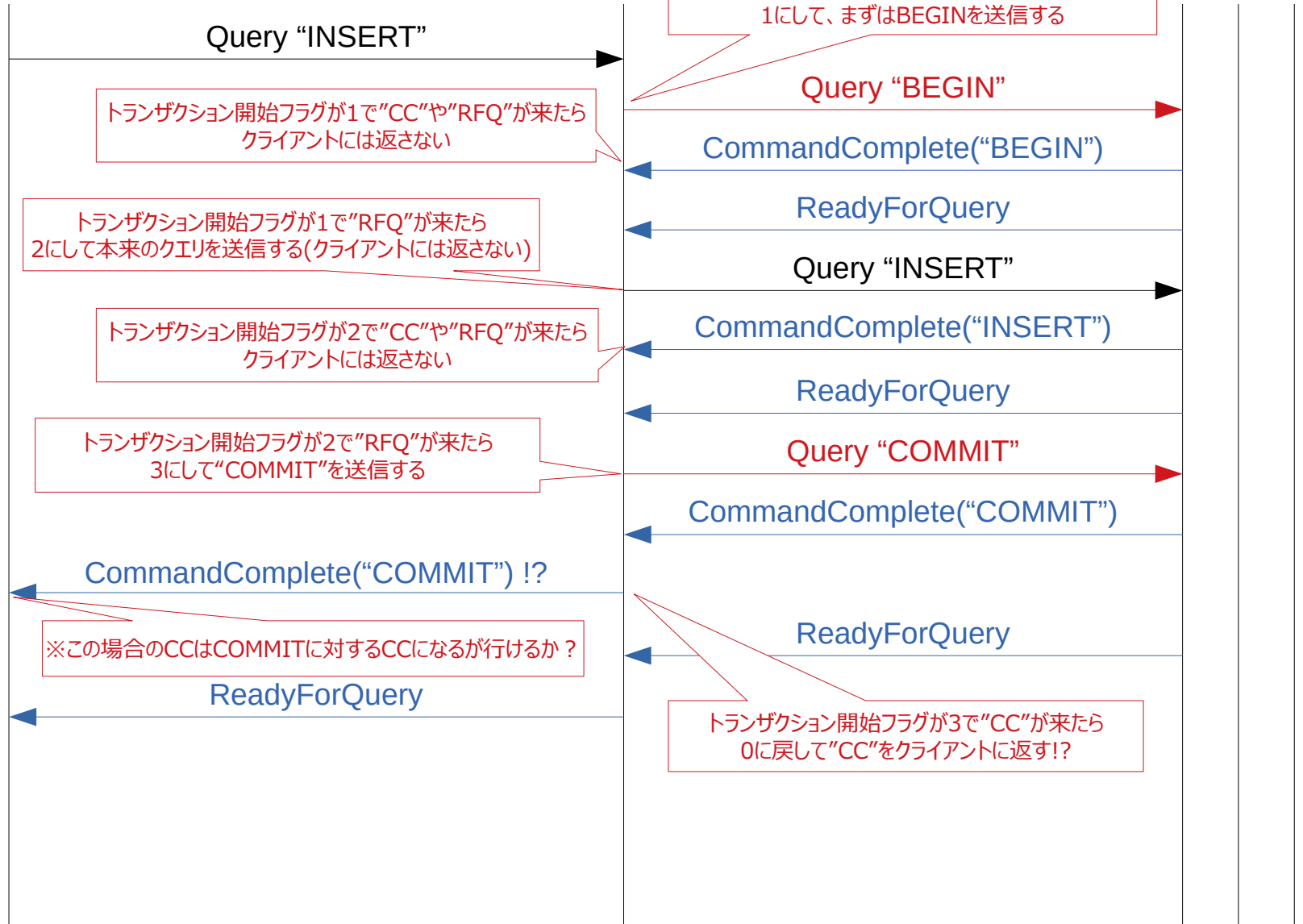
# PostgreSQL Multi DB Sequence Chart

PostgreSQL

Client

PgMulti

A B C



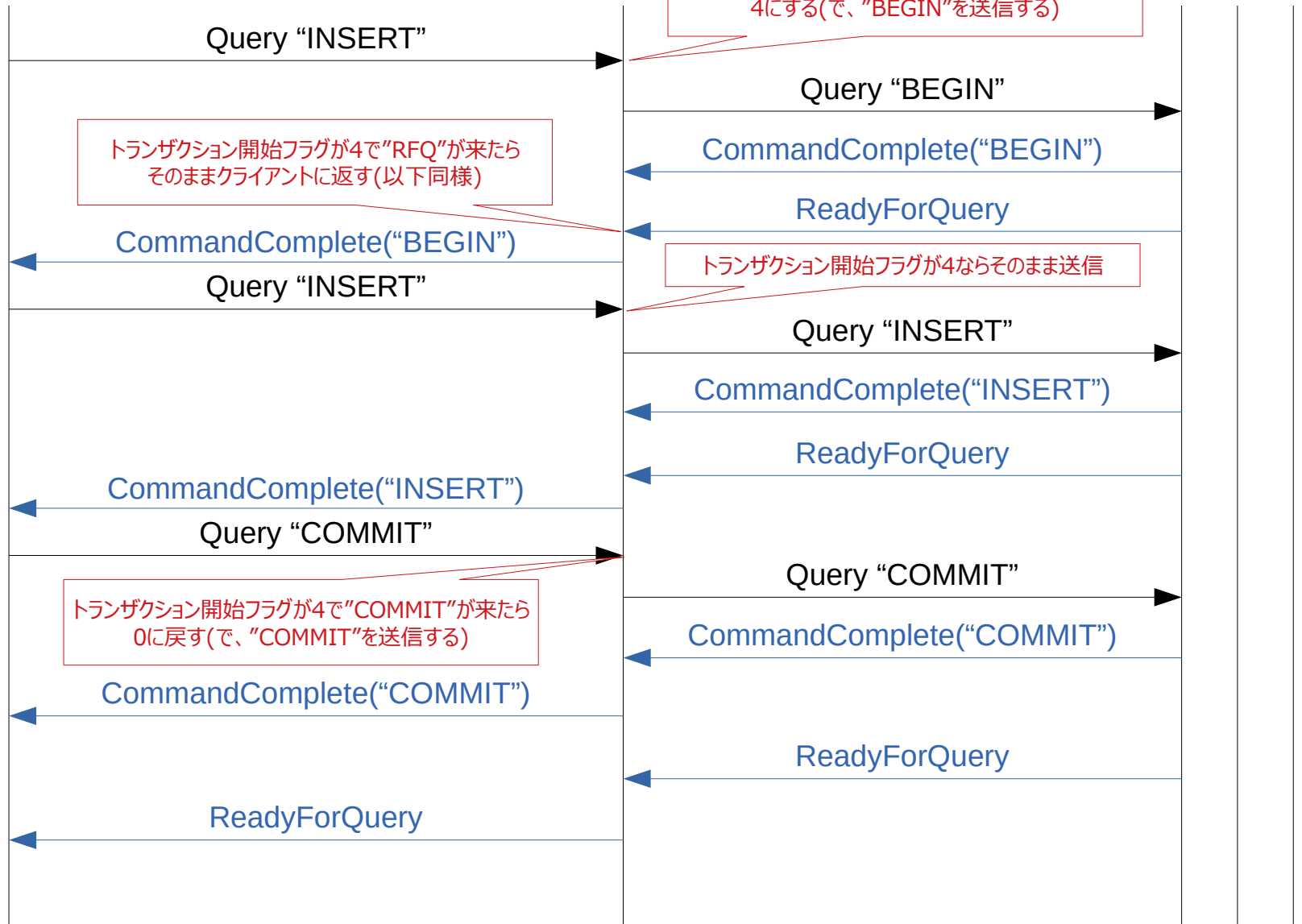
# PostgreSQL Multi DB Sequence Chart

PostgreSQL

Client

PgMulti

A B C



# ▶マルチメインDB目指して

- ・PgAnalyzerはGitHubに公開予定
- ・引き続き、まったりと本題の開発を進める予定

こんな感じで、PostgreSQL用のマルチメインDB  
実現用のミドルウェアを模索します。

たぶん機能とかを割り切って実装すれば、  
そんなに難しくなく実現できるはず…かなー？



ご期待ください!?





# おしまい

**Future Versatile Group/MyDNS.JP**  
**T.Kabu**

