

pgpool-II 2.0について

SRA OSS, Inc. 日本支社 石井 達夫

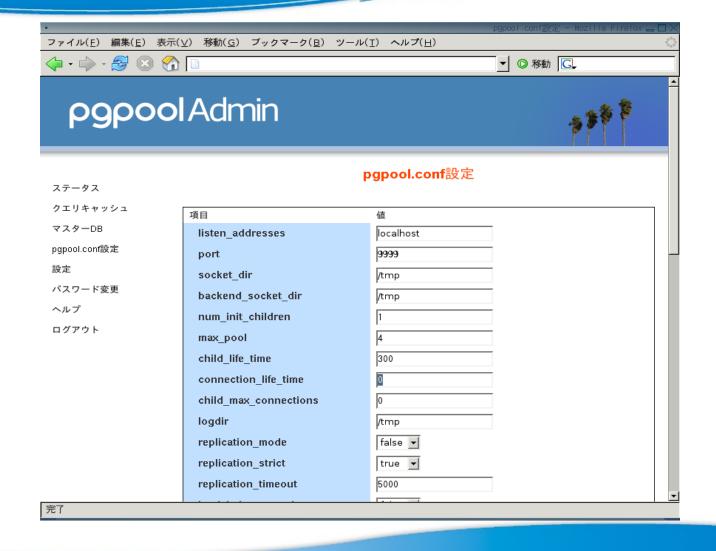
pgpool-IIとは?



- ◆pgpool Global Development Groupが開発
- ◆BSDライセンスで配布される
- ◆多彩な機能
 - ◆コネクションプーリング、レプリケーション、負荷分散、 パラレルクエリ
- ◆設定が容易. GUI管理ツールも附属
- ◆PostgreSQLに手を入れていないので、バージョン 追従が容易
- ◆使用言語を選ばない

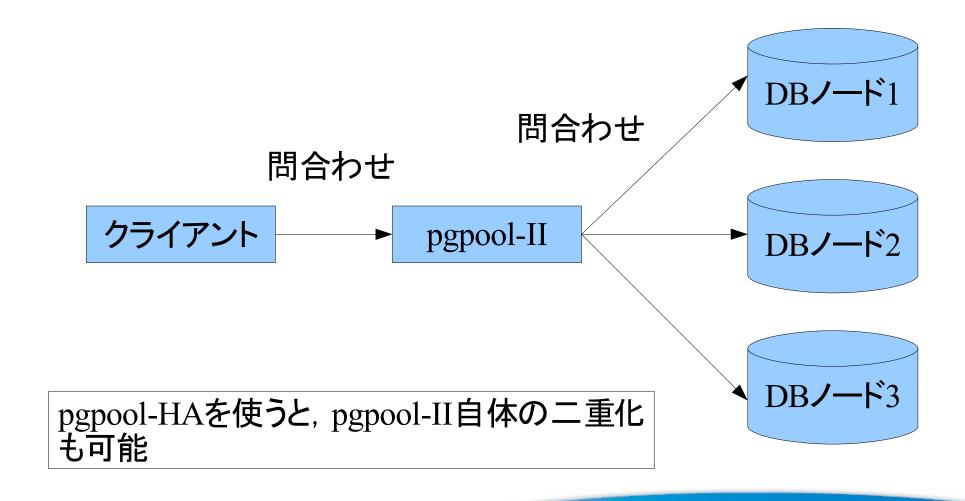


GUI管理ツール: pgpoolAdmin



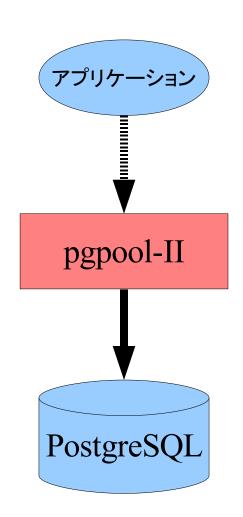
pgpool-IIのアーキテクチャ





pgpool-IIの機能

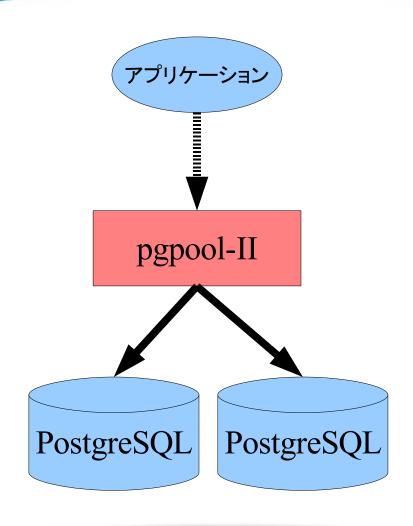
- ◆コネクションプーリング
 - ◆データベースへの接続 オーバヘッドを軽減
 - ◆PostgreSQLへのコネク ションを保持しておき、 再利用する
 - ◆Webシステムのように、 頻繁に接続/切断を繰り 返すシステムで効果あり



pgpool-IIの機能

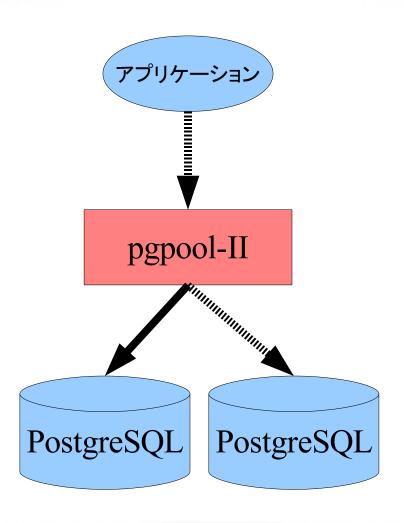
◆レプリケーション

- ◆SQL文を複製してデータベー スのコピーをリアルタイムに 作る
- ◆片方のDBが障害を起しても 自動的に切り離して運用を継 続
- ◆障害復旧後は運用を停めず にDBを同期,復帰させること が可能(オンラインリカバリ)
 - ◆オンラインリカバリの手 法をユーザが選択可能 なので、PITRを利用して 高速リカバリ可能



pgpool-IIの機能

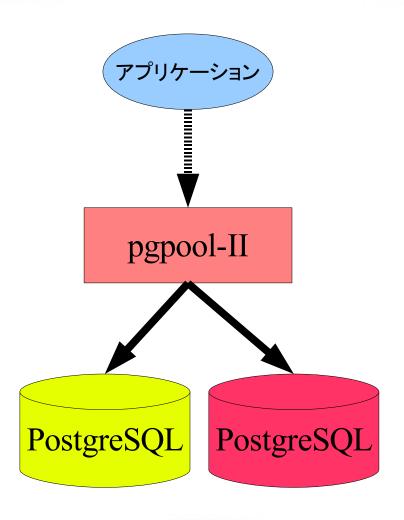
- ◆負荷分散
 - ◆検索問合わせをランダ ムに決めたPostgreSQL に振り向ける
 - ◆負荷をPostgreSQL間で 分かち合うので、検索性 能が向上
 - ◆TPC-Cなどの, OLTP(On Line Transaction Processing) 系の負荷に向いている



pgpool-IIの機能

◆パラレルクエリ

- ◆データを分割してPostgreSQL 間で分担
- ◆検索問合わせを複数の PostgreSQLで一斉に実行, 結果をpgpool-IIでまとめる
- ◆問合わせを並列に処理する ので高速
- ◆TPC-H/DBT-3などOLAP(On Line Analitical Processing)系の負荷に最適

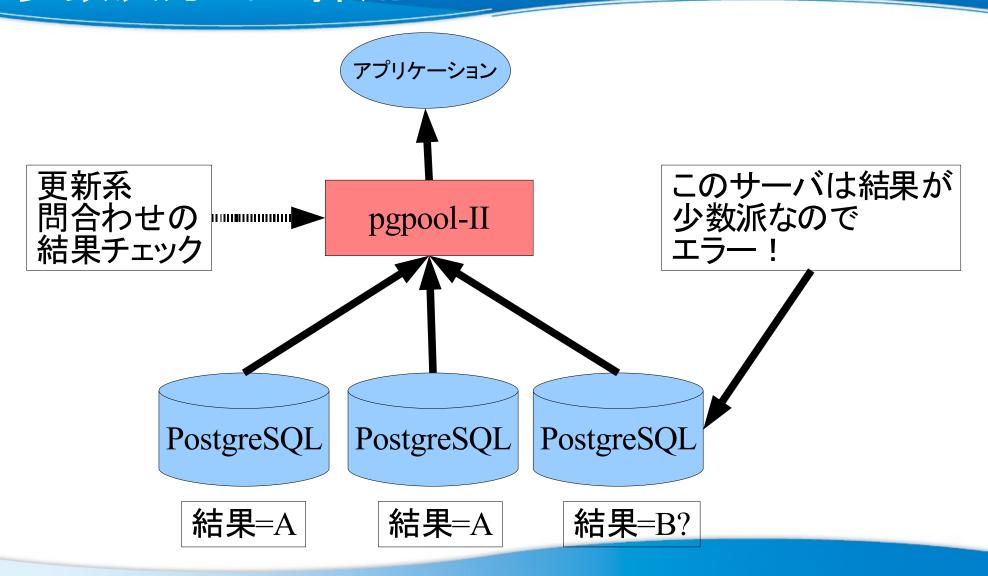


pgppool-II 2.0

- ◆2007/11/16 2.0 リリース
 - ◆その後マイナーバージョンアップして2.0.1がリリース済
 - http://pgfoundry.org/frs/download.php/1521/pgpool-II-2.0.1.tar.gz
- ◆様々な改良を追加
 - ◆レプリケーション
 - ◆信頼性向上、高速化、オンラインリカバリ
 - パラレルクエリ
 - ◆高速化
 - 全般
 - ◆他のシステムとの連携

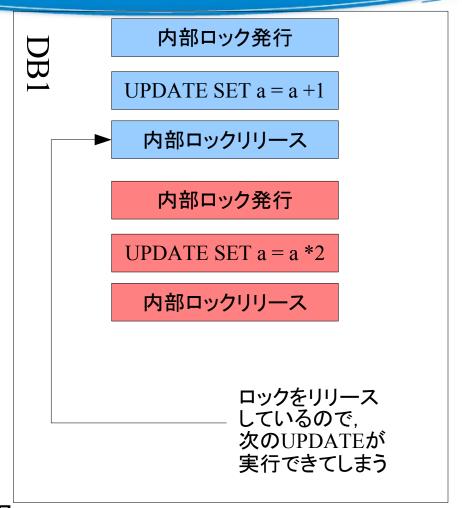
レプリケーションの信頼性向上: 多数決方式の採用





レプリケーションの信頼性向上: 明示的なトランザクション発行__





内部ロック発行 UPDATE SET a = a *2内部ロックリリース 内部ロック発行 UPDATE SET a = a + 1内部ロックリリース

時間

ノード間で実行順が違う!



レプリケーションの信頼性向上: 明示的なトランザクション発行



トランザクション開始 DB1 UPDATE SET a = a + 1DB2の処理が終わる までロックを保持 トランザクション終了 トランザクション開始 UPDATE SET a = a *2トランザクション終了 時間

トランザクション開始 UPDATE SET a = a + 1トランザクション終了 トランザクション開始 UPDATE SET a = a *2トランザクション終了

明示的なトランザクションを使うことでノード間の実行順を保証

レプリケーションの信頼性向上: 更新行数の確認



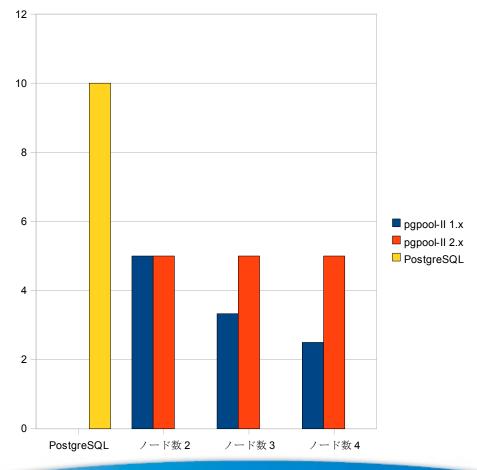
- ◆pgpool-II 1.xでは更新の成功, 失敗のみを判断
- ◆pgpool-II 2.0では、加えて更新件数もチェック
 - ◆件数が一致していない場合はエラー扱い
 - ◆回復はオンラインリカバリで

レプリケーションの高速化



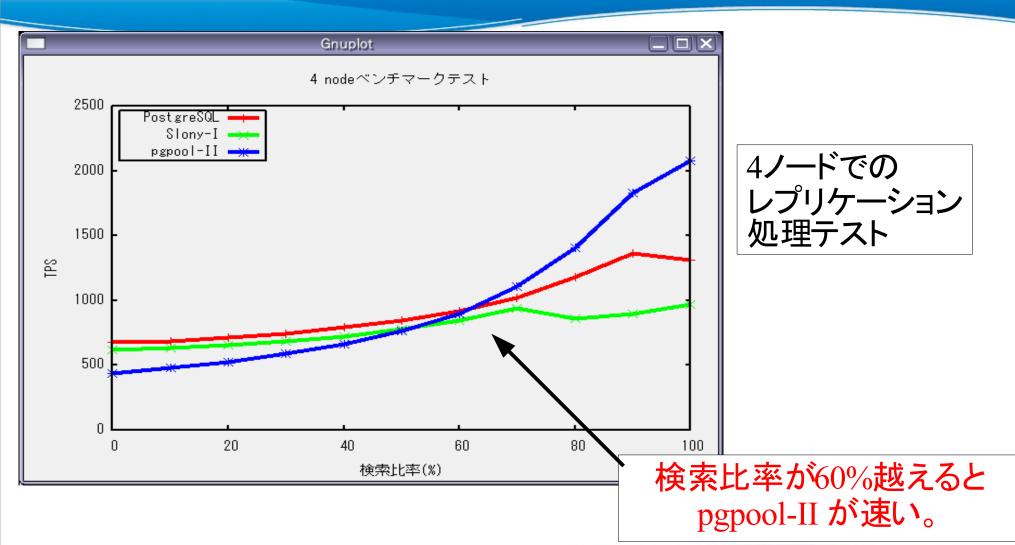
- ◆pgpool-II 1.xでは、更新性能はノード数に比例して低下
 - ◆ノード数が2なら1/2,3なら1/3,4なら1/4...
- ◆pgpool-II 2.0では, 更新性能はノード数に関わらず1/2

更新性能の比較



pgbench によるベンチマーク結果





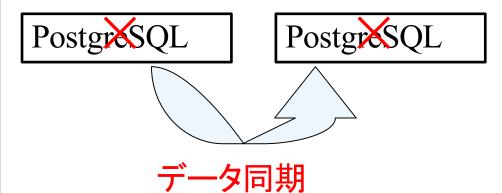


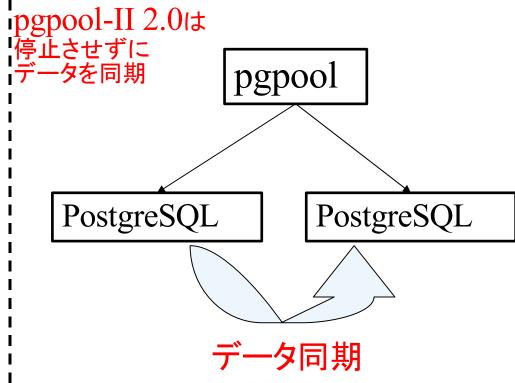
オンラインリカバリとは?

◆pgpoolを停止させずにデータを同期させてノードを復帰させる機能

pgpool-II 1.xは すべてを停止させて データを同期

pgpool





オンラインリカバリの仕組み



- ◆2段階に分けてデータを同期
 - ◆ファーストステージ
 - ◆並行して DB の参照・更新が可能
 - ◆ セカンドステージ
 - ◆すべての接続が終了するまで待機
 - 接続リクエストをすべてブロック
- ◆リカバリ方法
 - ◆ユーザ自身でどのようにリカバリするかを決めることができる
 - ◆pgpool-IIのソースコードにサンプルスクリプトを同梱

オンラインリカバリの実行例



- ◆ PITRを使ったオンラインリカバリ
 - ◆差分バックアップにより、接続をブロックする時間を短時間にすることが可能
 - ◆ブロックする時間はpostmasterがアーカイブログからリカバリする時間
- ◆ファーストステージ
 - ベースバックアップを取得し、ダウンしたホストへコピー
- ◆ セカンドステージ
 - ◆ SELECT pg_switch_xlog() を実行し、最新のWALログをアーカイブさせる

オンラインリカバリ設定



- pcp.conf
 - 🔷 pcp_recovery_node コマンドを使うので設定は必須
- 🔷 pgpool.conf
 - recovery_user
 - ◆リカバリ中にPostgreSQL(template1データベース)へ接続するためのユーザ名
 - recovery_password
 - ◆recovery_userのパスワード
 - recovery_1st_stage_command
 - ◆ファーストステージで実行するコマンド
 - ◆コマンドは\$PGDATA以下におく必要がある(セキュリティ上の問題)
 - recovery_2nd_stage_command
 - ◆セカンドステージで実行するコマンド
 - ◆コマンドは\$PGDATA以下におく必要がある(セキュリティ上の問題)
 - ◆ backend_data_directoryN(Nは数字)
 - ◆データベースクラスタディレクトリ名

オンラインリカバリ実行方法



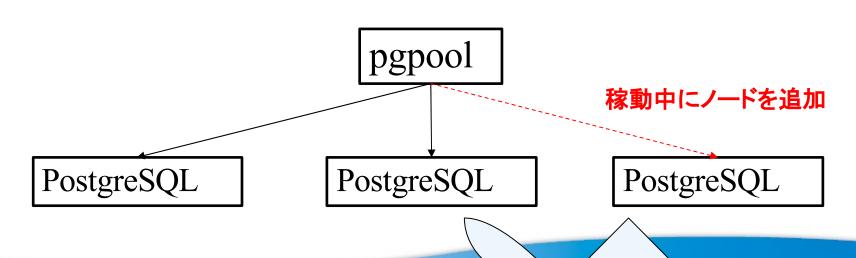
◆管理ツール上から実施

			pgpoolステータ	フス - Mozilla F	irefox 🕳 🗆 🗙	
ファイル(<u>F</u>) 編集(<u>E</u>)	表示(<u>V</u>) 移動(<u>G</u>)	ブックマーク(<u>B</u>)	ツール(<u>T</u>)	ヘルプ(<u>H</u>)	(2)	
pgpool Administration Tool						
▶ pgpoolステータス	pgpoolステータス	· ·			? ヘルプ	
▶ ノードステータス	サマリー プロセス情報 ノード情報					
▶ クエリキャッシュ	ノード情報					
▶ 分散ルール	IPアドレス オ	ポート ステータス	ウェ	イト		
▶ pgpool.conf設定	5-	432 ノード稼働中。	接続有り 0.33	3 切断		
▶ 管理ツール設定	5-	433 ノードダウン	0.33	3 リカ/	(i)	
トパスワード変更	5.	434 ノード稼働中。	接続無し 0.33	切断		
▶ ログアウト	サマリー ブ	プロセス情報 ノード情	 			
	pgpool pgpool再起動 設定リロード					

PowerGres 🕬

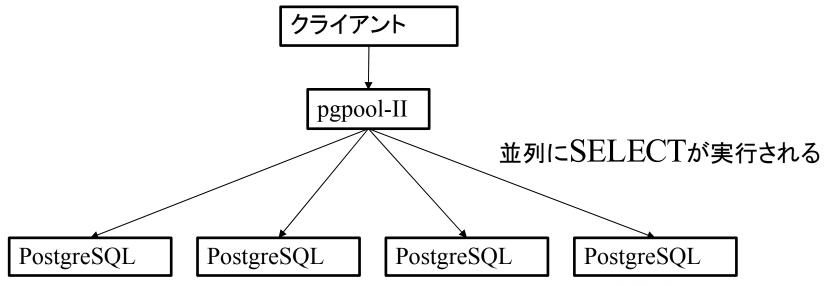
オンラインリカバリの応用

- ◆ノードの動的追加
 - 2.0 から pgpool.conf のリロードが可能
 - ◆backend_* を追加し、オンラインカバリすることで負荷状況に応じてノードを追加ができる(スケールアップ)
 - ◆リロードするとノードは「ダウン」状態で追加される



パラレルクエリとは?

- ◆複数のノードでリクエスト(SELECT)を処理
 - ◆1つのテーブルを複数サーバに分割(テーブルパー ティショニング)しておき、SELECTの結果をpgpool-II でまとめることで、仮想的に1つのテーブルとして扱う

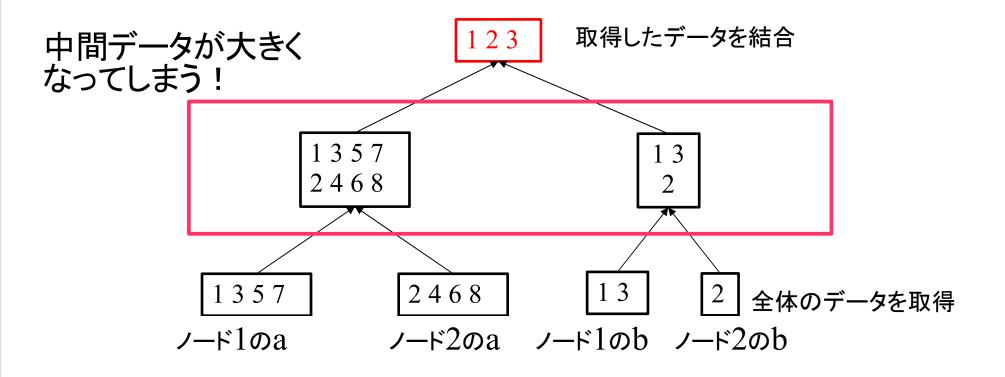


パラレルクエリの変更点

- ◆1.xの挙動
 - ◆すべてのテーブルを分割する必要がある
 - ◆中間データが大きくなるのでJOINが遅い
- **2.0**
 - ◆レプリケーションテーブルとパーティショニングテーブルを混在させることが可能
 - ◆クエリをより並列に動かすことが可能

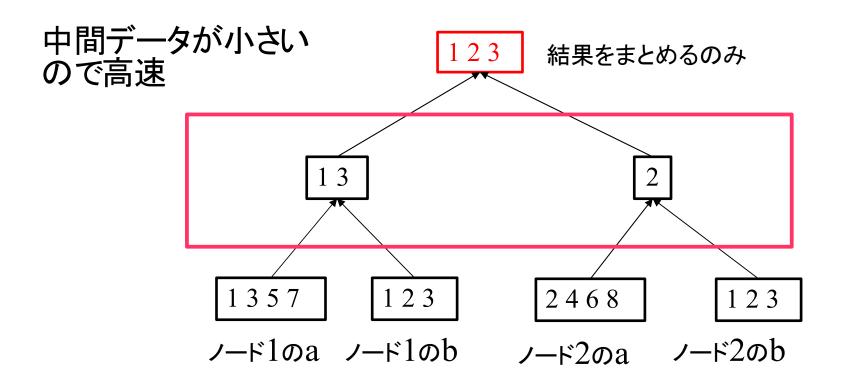
pgpool-II 1.xの動作







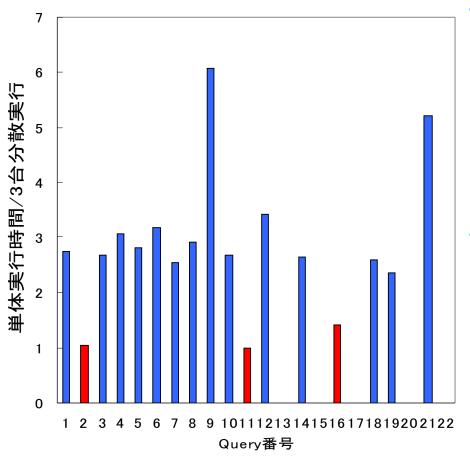
pgpool-II 2.0の動作



DBT-3によるベンチマーク



PostgreSQLと3台並列処理の実行時間の比較



- 🔷 line_itemとordersをテーブル分割
 - PostgreSQL 3 ノード(3分割)
 - ◆ 青線: パーティションテーブルとレプリケー ションテーブルの結合
 - ◆ 赤線:レプリケーションテーブル同士の結合
 - ◆ パラレルクエリの弱いクエリは未測定
 (13.15.17.20.22)
- scale factor=10
 - ◆インデックス含めてトータルで30GB

最大で6倍のパフォーマンス改善 (クエリ9番)

他のシステムとの連携

- failover_command
 - ◆ノードを切り離した時に指定したコマンドを実行
 - ◆Slony-Iの連携やアラートメールを送ることが可能
- failback_command
 - ◆ノードを戻した時に指定したコマンドを実行
- ◆ pgpool-IIから必要な情報を渡すことも可能(特殊文字を置換)
 - ◆%d:対象ノード番号
 - ◆%h:対象ノードのホスト名
 - ◆%p:対象ノードのポート番号
 - ◆ %D:対象ノードのデータベースクラスタパス

クライアントの強制切断

- client_idle_limit
 - ◆クライアントからクエリが届くまでの最大待ち時間を 指定可能
 - ◆秒単位で指定
 - ◆TCPのタイムアウト待ち防止などに有効

pgpoolファミリの今後



- ◆pgpool-Iは今後pgpool-IIに移行してもらう
- ◆pgpool-II 1.xは今後メンテを終了予定
- ◆pgpool-II 2.xへ一本化
- ◆pgpool-HAはpgpool-IIでも利用可能



- ◆pgpool-II 2.0の特長
 - ◆レプリケーションの信頼性向上
 - ◆更新処理の高速化(4ノードなら従来の2倍!)
 - ◆オンラインリカバリ・ノード追加
 - ◆部分レプリケーションによるパラレルクエリの高速化
- ◆SRA OSSでは、pgpool-IIの商用サポートを提供する予定です!

参考URL



- ◆ pgpool-II 開発サイト
 - http://pgfoundry.org/projects/pgpool/
- pgpool Wiki
 - http://pgpool.sraoss.jp/
- ◆ pgpoolメーリングリスト
 - http://www.sraoss.jp/mailman/listinfo/pgpool-general-jp
- ◆ マイコミジャーナルでのpgpool-IIの解説記事
 - http://journal.mycom.co.jp/articles/2007/11/08/pgpool/index.html

pgpool-II デモシステムの構成



