

PGCluster Round Table

PGCluster円卓会議

At.Mitani

PostgreSQL 10th Anniversary Summit

9 July 2006 Toronto

Agenda

- エンタープライズ規模で利用可能にしてほしいという要望がある
- レプリケーションがそんな要求をみたすことができるのか？ Requirement?
- PGClusterについて
 - Feature 機能
 - Structure レプリケーションの仕組み
 - demonstration デモ
 - Replication レプリケーションの将来
 - HA feature 高可用性機能
- PGClusterへのその他の要望 Requirement for PGCluster
- PGClusterIIのコンセプト PGCluster-II
 - Feature 機能
 - How to 分散サーバ間の共有IPCの方法 distributed servers
 - Structure 構造
 - demonstration デモ
 - Expectation effect and problem 期待される効果と問題点
- 質疑応答 Q&A Session

Issues of Enterprise usage

エンタープライズ規模での利用上の問題

Requirement of Enterprise usage

- 巨大なデータの取り扱い
— SANやNASを利用した巨大なストレージの使用
- 大量のリクエストに対する迅速なレスポンス
— 検索（例：ウェブアプリ）
— 更新（例：ERPシステム）
- 高可用性
— 無停止オペレーション
— 重大な事故からの復旧
- 簡単なバックアップとリストア
— 差分バックアップ/リストア

Does replication help of those requirement ! (1/4)

レプリケーションはこんなことの役に立つのか？ (1/4)

- 巨大なデータの操作 data operation

- あまり意味がない

- 操作するデータの量を減らすものではありません amount of operation data

- その他のソリューション on

- PostgresForest / pgpool-II

Does replication help of those requirement ? (2/4)

レプリケーションはこんなことの役に立つのか？ (2/4)

- 負荷分散 | distribution

- 効果があるところ

- 検索の負荷が軽減します search load

- あまり効果がないところ

- 更新の負荷はもっと上がってしまいます！ update heavy!

- その他のソリューション on

- クラスタ/グリッドコンピューティング
- これで鉄板というツールはわかりません other tools

Does replication help of those requirement ? (3/4)

レプリケーションはこんなことの役に立つのか？ (3/4)

- 高可用性 high availability

- 効果あり

- マスタ-スレーブ構成では、テイクオーバーの際にサービスが停止します service stop time during take over.
- 複数マスタ構成では、テイクオーバーの際にもサービス停止は発生しません service stop time during take over.

- その他のソリューション on

- サーバレベルの高可用性のための商用ツールがあります for server level HA.

Does replication help of those requirement ? (4/4)

レプリケーションはこんなことの役に立つのか？ (4/4)

- 簡単なバックアップとリストア and restore
 - 効果あり
 - レプリケーションされたデータベースはリアルタイムのバックアップになります
 - その他のソリューション n
 - 差分バックアップとリストア up / restore

PGCluster / PGCluster-II

PGClusterとPGCluster-II

History 歴史

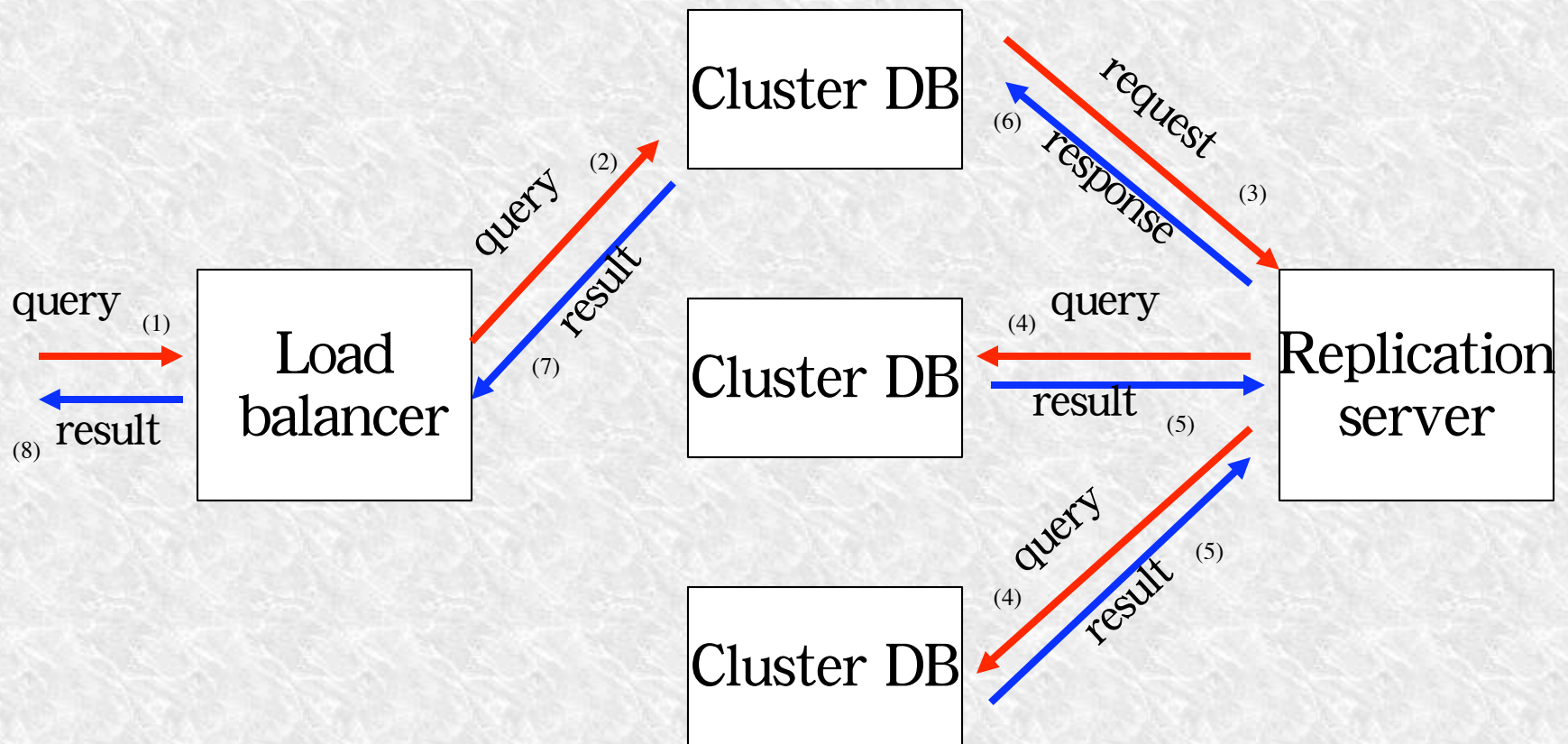
- PGReplicate-0.1 released (2002 Mar)
 - クエリベースの単純なレプリケーション query
- PGCluster-1.0 released (2003 July –)
 - トランザクション、コピー、ストアードプロシージャ lure
- PGCluster-1.1 released (2004 Sep –)
 - ロードバランサにPgpoolを利用 lancer
- PGCluster-1.3 released (2004 Dec –)
 - ラージオブジェクト
- PGCluster-1.5 released (2006 Jan –)
 - V3 プロトコル、プリペアドクエリ ury
- PGCluster-II started (2006 June –)

Features of PGCluster

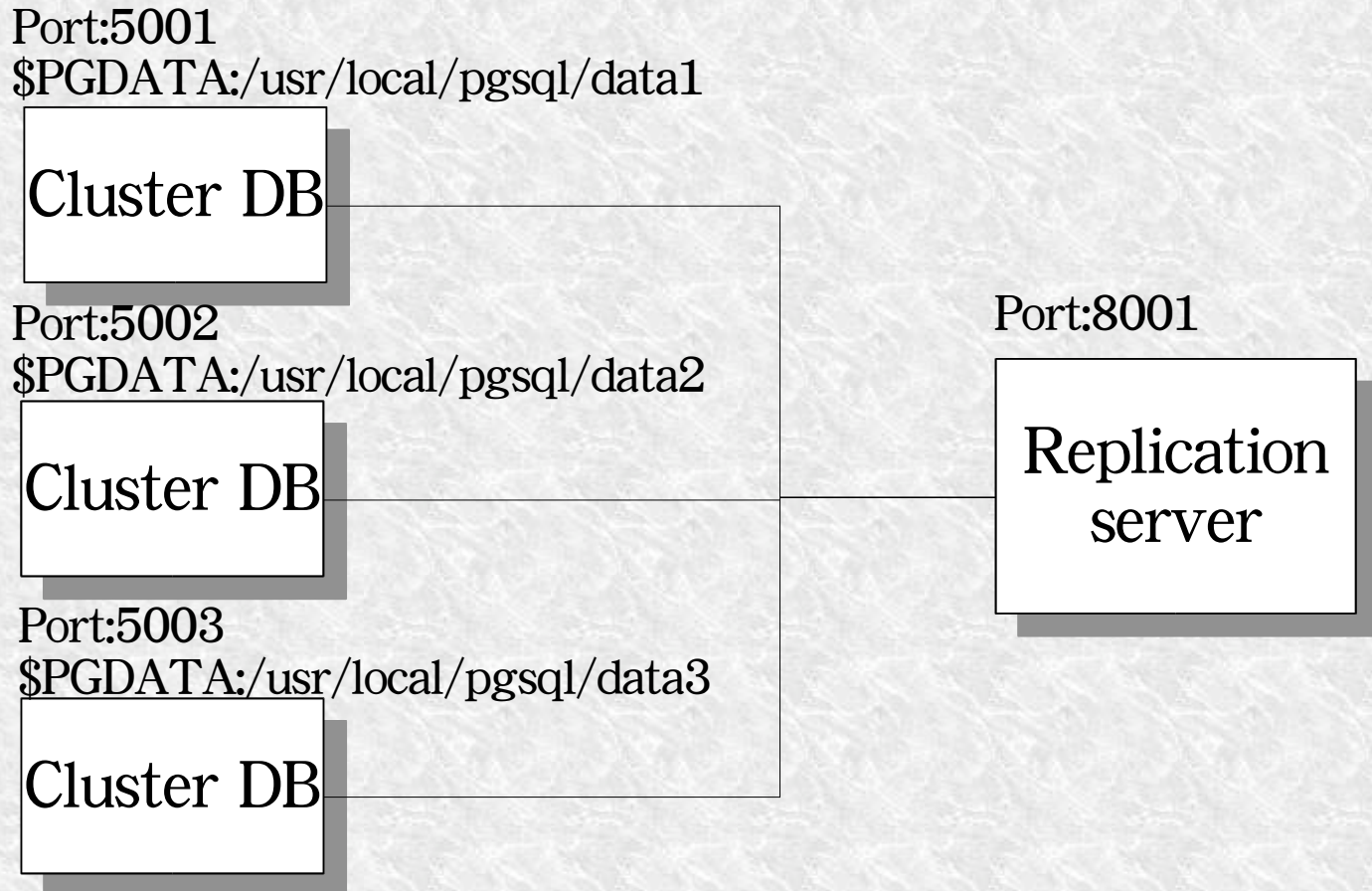
PGClusterの機能

- マルチマスタ & 同期レプリケーションシステム
synchronous replication system
 - 適している用途 for
 - 読み込みの負荷が高いシステム
 - 高可用性の求められるシステム
 - リアルタイムバックアップ
 - 適していない用途 not suitable for
 - 書き込みの負荷が高いシステム
 - 巨大なサイズのデータベース
 - Mobile database 構成変更の多いデータベース

Structure of Replication

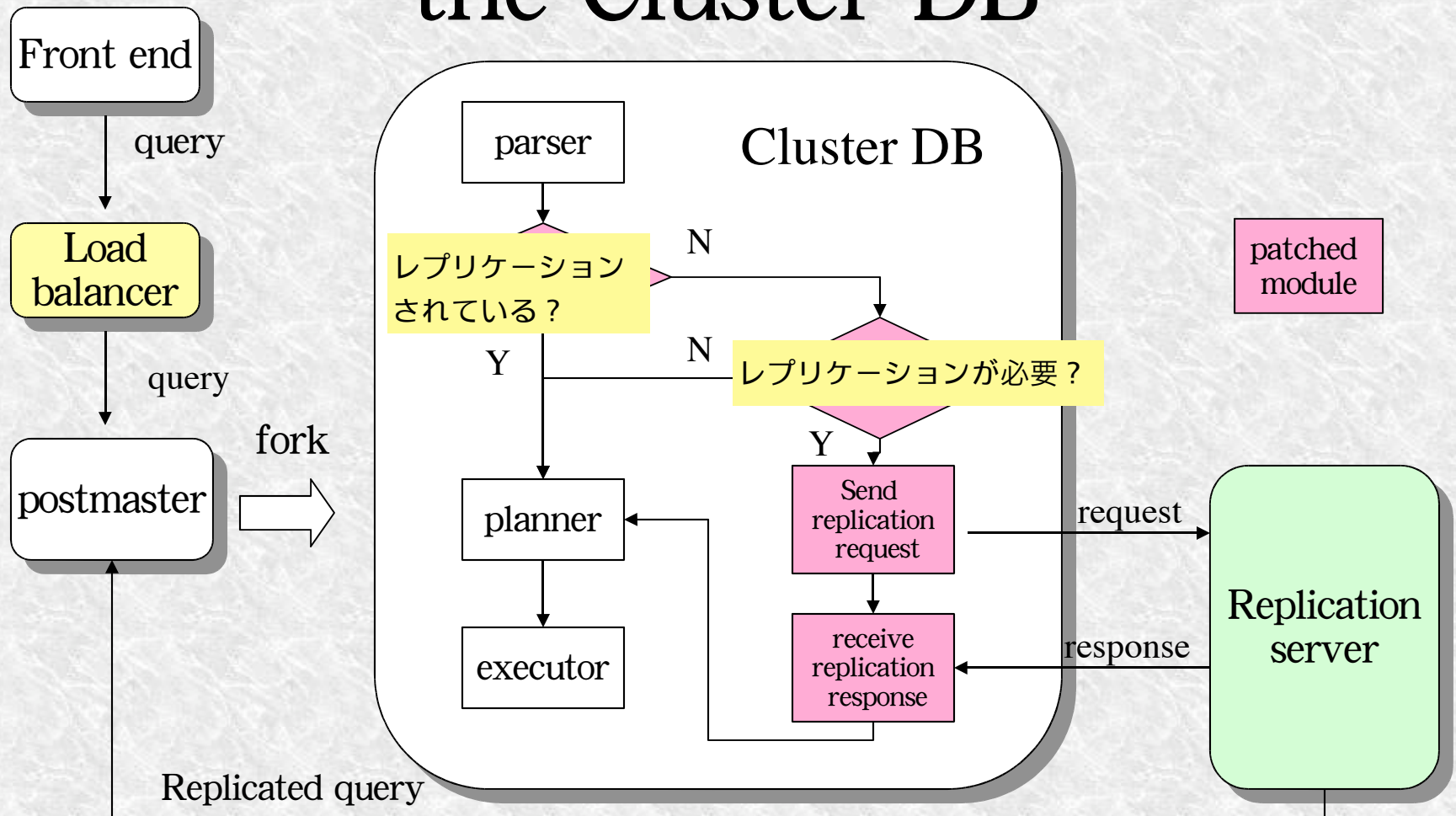


Demonstration of PGCluster



Replication patch of the Cluster DB

クラスタDBでレプリケーションが担当するところ



What cannot be replicated

レプリケーションされるもの

- トランザクション
- ストアドプロシージャ
- 内部関数の戻り値 value of internal function
 - now(), random()
- ラージオブジェクト
- プリペアドクエリ
- Copy (from)
- Create / Drop
 - DB , table, domain, function, group ...

High Availability support

高可用性のサポート

- 自動テイクオーバー c take over
 - クラスタDB DB
 - レプリケーションサーバ erver
- 動的なサーバ追加 ; server addition
 - クラスタDB DB
 - レプリケーションサーバ erver
- 代替のない破損箇所がないこと ngle point of failure
 - 全てのサーバで複数台構成が可能 et up multiple

Yet another requirement for PGCluster

さらにPGClusterへの要望

- 書き込みの負荷分散 dispersion
- 巨大なサイズのデータベース atabase
- リカバリの時間短縮 recovery

Write load dispersion & Large size of database

書き込みの負荷分散と巨大なサイズのデータベース

- レプリケーションは書き込みの負荷分散の役には立たない for write load distribution
 - 現在のバージョンのスコープにはない for current version.
- 必要とされる分散データベース distributed database
 - パラレル操作 parallel operation
 - PostgresForest
 - Pgpool-II
 - 分散サーバ間の共有IPC during distributed servers

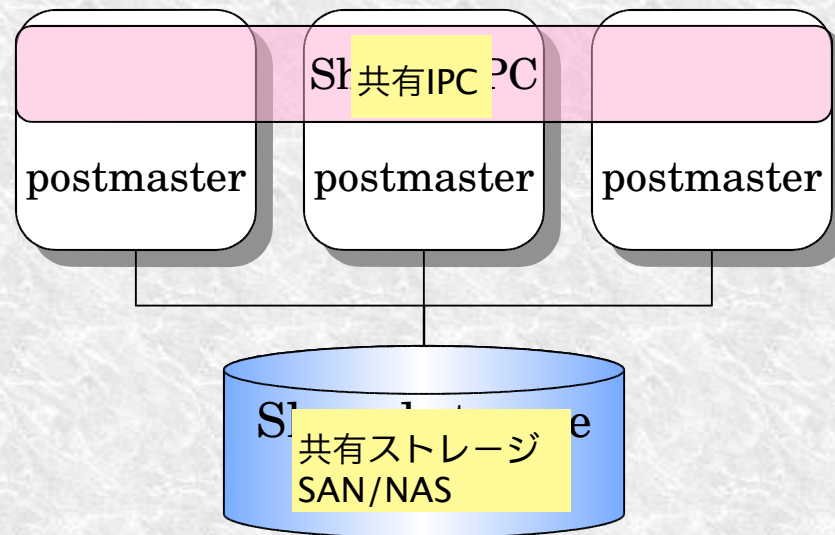
Short recovery

リカバリの時間短縮

- システム変更の多いユーザからの要求 file users
 - The files under \$PGDATA are copied using
 - rsync command
 - 差分のあるファイルのみコピーする a difference is copied
 - 差分が小さくても、ファイル全体がコピーされる, a file is all copied.
 - アーカイブログを利用すべき... should be used...

Concept PGCluster-IIのコンセプト cluster-II

- 共有DBクラスタ DB cluster
 - 共有ストレージ (SAN/NASを利用) using SAN / NAS)
- 分散postmaster distributed postmaster
 - 分散サーバ間の共有IPC during distributed servers

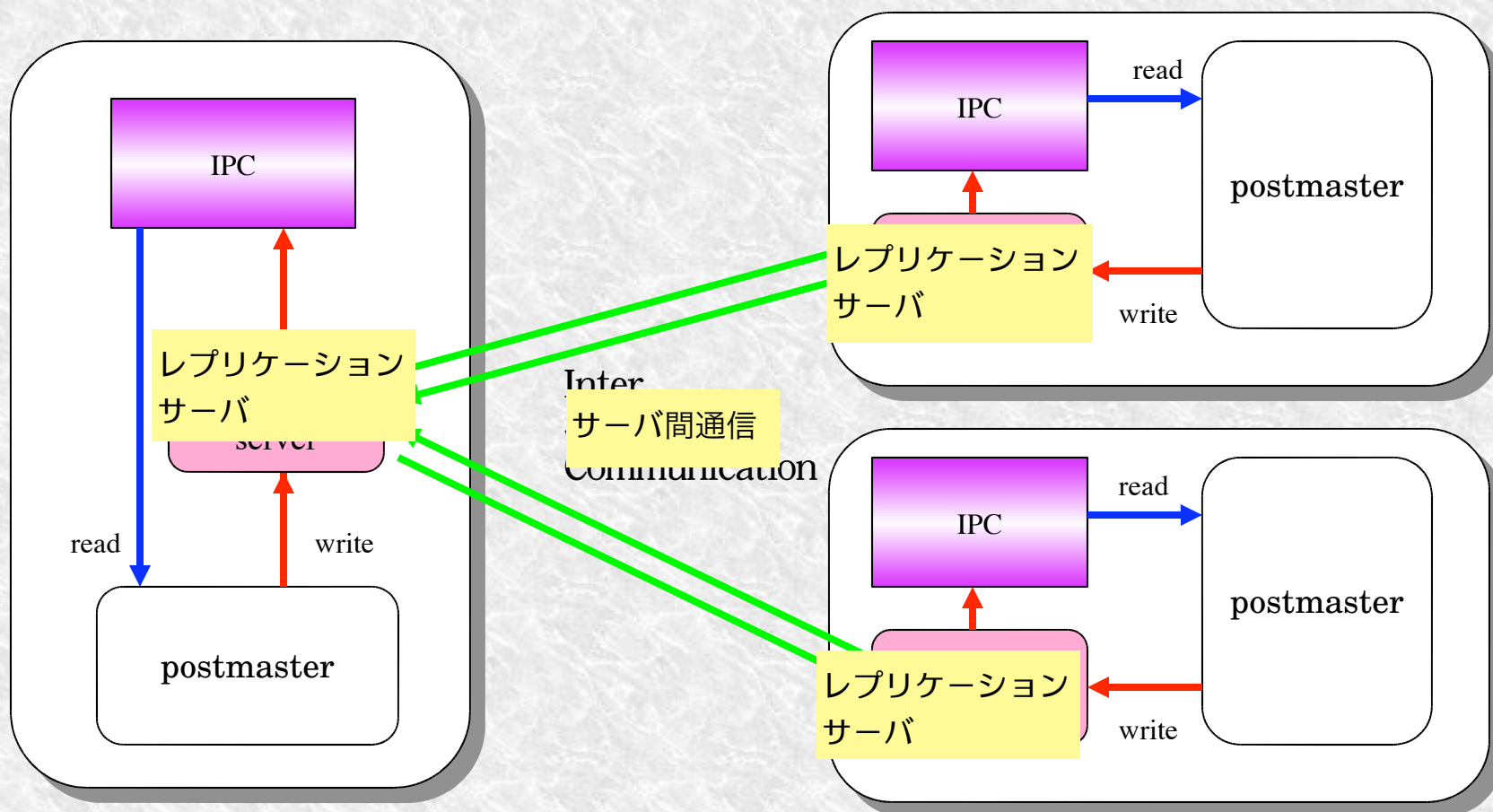


How to share IPC during distributed servers

分散サーバ間でどのようにIPCを共有するか

- 基本的にIPCはプロセス間通信の方法としては最速 the fastest communication method during process
 - ゆえに、どのような方法で共有してもそれがオーバーヘッドとなる becomes an overhead
- 可能な限りローカルのIPCを利用する local IPC as much as possible
 - 書き込みの同期 (レプリケーション) only (replication)
 - マルチマスタ & 同期レプリケーションが必要 master & synchronous replication
 - ローカルで読み出し only

Server composition of PGCluster-II



Patch point in order to share semaphore

セマフォ共有のためのパッチポイント

| Interface function | Source | Target function |
|-----------------------|-----------------------|-----------------------------------|
| PGReserveSemaphores() | storage/ipc/ipci.c | CreateSharedMemoryAndSemaphores() |
| PGSemaphoreCreate() | storage/lmgr/spin.c | SpinlockSemas() |
| PGSemaphoreReset() | storage/lmgr/proc.c | ProcCancelWaitForSignal() |
| | | LockWaitCancel() |
| PGSemaphoreLock() | storage/lmgr/lwlock.c | LWLockAcquire() |
| | storage/lmgr/proc.c | ProcSleep() |
| | | ProcWaitForSignal() |
| PGSemaphoreUnlock() | storage/lmgr/lwlock.c | LWLockAcquire() |
| | | LWLockRelease() |
| | storage/lmgr/proc.c | ProcWakeup() |
| | | CheckDeadLock() |
| | | ProcSendSignal() |
| | storage/lmgr/spin.c | s_unlock_sema() |
| PGSemaphoreTryLock() | storage/lmgr/spin.c | tas_sema() |

6 delegate functions

Called 14 functions

Patch point in order to メモリ共有のためのパッチポイント share memory

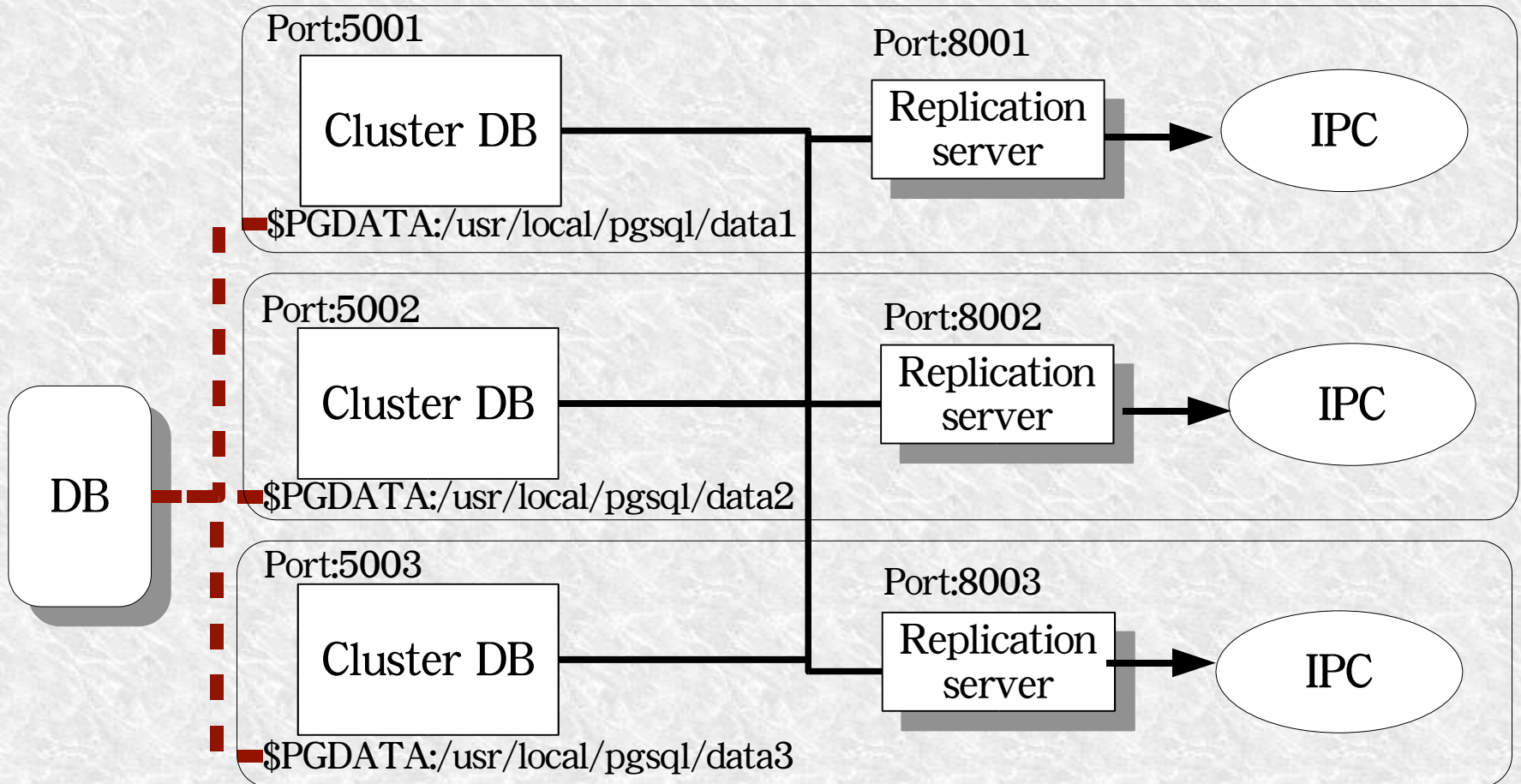
| | | |
|----------------|----------------------------|--|
| BgWriterShmem | postmaster/bgwriter.c | BackgroundWriterMain() BgWriterShmemInit() |
| MultiXactState | access/transam/multixact.c | GetNewMultiXactId() MultiXactShmemInit() StartupMultiXact() MultiXactSetNextMXact() MultiXactAdvanceNextMXact() TruncateMultiXact() |
| TwoPhaseState | access/transam/twophase.c | TwoPhaseShmemInit() MarkAsPreparing() RemoveGXact() |
| ControlFile | access/transam/xlog.c | XLogWrite() WriteControlFile() ReadControlFile() UpdateControlFile() XLOGShmemInit() BootstrapXLOG() StartupXLOG() CreateCheckpoint() |
| XLogCtl | access/transam/xlog.c | XLogInsert() AdvanceXLogInsertBuffer() XLogWrite() XLOGShmemInit() StartupXLOG() |

| | | |
|-------------------|---|---|
| BufferDescriptors | storage/buffer/buf_init.c storage/buffer/bufmgr.c storage/buffer/freelist.c | InitBufferPool() ReadBuffer() PinBuffer() StrategyGetBuffer() StrategyFreeBuffer() |
| BufferBlocks | storage/buffer/buf_init.c | |
| StrategyControl | storage/buffer/freelist.c | StrategyGetBuffer() StrategyFreeBuffer() StrategyInitialize() |
| FreeSpaceMap | storage/freespace/freespace.c | InitFreeSpaceMap() delete_fsm_rel() realloc_fsm_rel() link_fsm_rel_usage() unlink_fsm_rel_usage() link_fsm_rel_storage() unlink_fsm_rel_storage() compact_fsm_storage() push_fsm_rels_after() |
| PMSignalFlags | storage/ipc/pmsignal.c | PMSignalInit() SendPostmasterSignal() CheckPostmasterSignal() |
| procArray | storage/ipc/procarray.c | CreateSharedProcArray() ProcArrayAdd() ProcArrayRemove() |

| | | |
|----------------|-------------------------|---|
| shminvalBuffer | storage/ipc/sinvaladt.c | SIBufferInit() |
| newLockMethod | storage/lmgr/lock.c | LockMethodTableInit() |
| ProcGlobal | storage/lmgr/proc.c | InitProcGlobal() ProcKill() DummyProcKill() |
| DummyProcs | storage/lmgr/proc.c | InitProcGlobal() |

15 pointers
Called from 50 functions

Demonstration of PGCluster-II



Expected effect and problem

期待される効果と問題点

- 効果 effect
 - 書き込みの負荷分散 decrease improvement of write load
- 問題点 blem
 - 共有IPCのオーバーヘッド increase the shared IPC
 - ハードウェアコストの上昇 increase hardware cost
 - ストレージ (NAS/SAN) / SAN)
 - ネットワーク (iSCSI/FC) / FC)

Discussion

質疑応答

Discussion (Enterprise usage)

質疑応答（エンタープライズ規模での利用）

- いかにしてパフォーマンスを改善するか the performance
 - データサイズが巨大な場合 the size of data
 - アクセスが多い場合 (search/update/insert)
- いかにして可用性を確保するか the availability
 -
- いかにして短時間で巨大なデータをバックアップすることができるか data in a short time
- どうやってこれら全てを実現するか all of them

Discussion (PGCluster)

- サードパーティとして貢献するには
— 3rd party
- 共有 / 非共有
— Shared / Non shared
 - どちらがより重要なのか
— Which is more important
- PGClusterに必要なのはなにか
— What is needed for PGCluster

Thank¹ you

謝辞

- When you have a question || issue || requirement for PGCluster, would you please send email.
 - mitani@sraw.co.jp
 - pgcluster-general@pgfoundry.org
- You can download all version of PGCluster from following site,
 - <http://pgfoundry.org/projects/pgcluster/>