

A Middleware to Synchronize Multiple Instances in a PostgreSQL Cluster

SRA OSS, Inc. Japan Takeshi MISHIMA



About me

- I had been researching database systems at NTT and Keio University in Japan.
- My representative papers were accepted by VLDB2009 and SIGMOD 2015.
- I received Ph.D. degree from The University of Tokyo.
- Now, I work for SRA OSS, Inc. Japan.







1. What is Pangea?



2. The disadvantage of streaming replication



3. The algorithms of Pangea





PART 1

What is Pangea?

What is Pangea?

• Proposed Pangea and accepted by VLDB2009.

A synchronization middleware for multiple Postgres instances without modifying the source codes of Postgres.

- Pangea is a just prototype but the algorithms have been implemented into the latest version of Pgpool-II with the snapshot isolation mode.
- In today's my presentation, I will explain Pangea, and enjoy the fruit by using Pgpool-II if you are interested in Pangea !



The advantages of Pangea

• The advantages

(1) not only it increases availability but also,

(2) any instance returns the latest consistent data.

cf. Streaming replication provides (1) but does not (2).

Next, I will explain the detail of it.





PART 2

The disadvantage of streaming replication























• Even if we have a lot of standby instances, all queries must be sent to only the primary to ensure consistency.









The goal of Pangea

• Even if Pangea sends a SELECT query to any instance, it must return the latest consistent data.







The goal of Pangea

• Even if Pangea sends a SELECT query to any instance, it must return the latest consistent data.







PART 3

Pangea's algorithms to address the disadvantage

Assumption for Pangea

- Snapshot isolation (REPEATABLE READ isolation level)
- A query does not have any non-deterministic functions such as RANDOM and NOW.
- We call one replica the leader and the others the followers.
- We call the first DML query a snapshot query.





The replication method of Pangea

• Instead of using streaming replication, Pangea sends an UPDATE query to all instances.







Recap1

A snapshot is taken just before a snapshot query (the first DML query) is executed.





Recap2

All queries are executed on the snapshot.





Recap2

Update queries are also executed on the snapshot.





Recap3

The updates are applied to the original database when the transaction commits.





- Pangea propagates snapshot and commit queries exclusively.
- That is, Pangea propagates commit queries only when all instances do not execute snapshot queries.
- Conversely, Pangea propagates snapshot queries only when all instances do not execute commit queries.





 When Pangea sends commit queries to all instances, it stops sending snapshot ones.







• When Pangea sends snapshot queries to all instances, it stops sending commit ones.









This algorithm introduces that all instances must provide the same snapshot data.







This algorithm introduces that all instances must provide the same snapshot data.







This algorithm introduces that all instances must provide the same snapshot data.







This algorithm introduces that all instances must provide the same snapshot data.







This algorithm introduces that all instances must provide the same snapshot data.







The exclusive execution of snapshot and commit queries makes all instances create the same snapshot!





• First, Pangea propagates all UPDATE queries to only the leader.



















• First, Pangea propagates all UPDATE queries only to the leader.







 Second, if Pangea receives the answer, it sends the query to all the followers.







- This protocol introduces that write/write conflicting queries are serialized by the leader,
- and non-conflicting queries are executed on all instances simultaneously.





In a nutshell

- Algorithm1 makes all instances create the same snapshot, and
- Algorithm2 makes the order of only conflicting queries serialized by the leader.

- Pangea's algorithms keep consistency between the leader and followers loosely.
- Any instance can return the latest consistent data.
- Surprisingly, Pangea's algorithms need not modify the source codes of Postgres!







PART 4 Summary

Summary

- We proposed the algorithms to address the disadvantage of streaming replication.
- Formally, we gave the full proof of the correctness in our VLDB2009 paper, please see it for detail.
- Our algorithms have been implemented into the latest version of Pgpool-II, and therefore, enjoy the fruit of Pangea by using Pgpool-II !







mishima@sraoss.co.jp



