

BDR for DBAs

Martín Marqués

2ndQuadrant

October 23, 2015

Content

- 1 Replication history with PostgreSQL
- 2 BDR
- 3 BDR join and part nodes
- 4 BDR Sequence
- 5 BDR prohibited DDL commands
- 6 Conclusions
- 7 Questions?

Trigger based Replication

Trigger based Replication

- Trigger based

Trigger based Replication

- Trigger based (Slony, Londiste)

Trigger based Replication

- Trigger based (Slony, Londiste)
 - Used since 7.3

Trigger based Replication

- Trigger based (Slony, Londiste)
 - Used since 7.3
 - Trigger based (double writes on upstream)

Trigger based Replication

- Trigger based (Slony, Londiste)
 - Used since 7.3
 - Trigger based (double writes on upstream)
 - Needs external daemon

Trigger based Replication

- Trigger based (Slony, Londiste)
 - Used since 7.3
 - Trigger based (double writes on upstream)
 - Needs external daemon
 - No DDL replication

Stream Replication

- Stream Replication

Stream Replication

- Stream Replication
 - Available since 9.0

Stream Replication

- Stream Replication
 - Available since 9.0
 - One way replication

Stream Replication

- Stream Replication
 - Available since 9.0
 - One way replication
 - Hot Standby with 9.1

Stream Replication

- Stream Replication
 - Available since 9.0
 - One way replication
 - Hot Standby with 9.1
 - Can't replicate per-DB

Replication Slots

- Logical decoding

Replication Slots

- Logical decoding
 - Available since 9.4

Replication Slots

- Logical decoding
 - Available since 9.4
 - Good option for trigger based replication tools

Replication Slots

- Logical decoding
 - Available since 9.4
 - Good option for trigger based replication tools
 - BDR, UDR

Content

- 1 Replication history with PostgreSQL
- 2 BDR
- 3 BDR join and part nodes
- 4 BDR Sequence
- 5 BDR prohibited DDL commands
- 6 Conclusions
- 7 Questions?

Bi-Directional Replication

Bi-Directional Replication

- Per-database replication

Bi-Directional Replication

- Per-database replication
- Replication sets

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins \Rightarrow user-defined conflict handlers

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins \Rightarrow user-defined conflict handlers
- DDL replication

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins \Rightarrow user-defined conflict handlers
- DDL replication
 - Keeps schema synced

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins \Rightarrow user-defined conflict handlers
- DDL replication
 - Keeps schema syncd
 - ... even objects not in replication sets

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins \Rightarrow user-defined conflict handlers
- DDL replication
 - Keeps schema synced
 - ... even objects not in replication sets
 - DDL locks

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins \Rightarrow user-defined conflict handlers
- DDL replication
 - Keeps schema synced
 - ... even objects not in replication sets
 - DDL locks
 - Global lock will cancel queries

Bi-Directional Replication

- Per-database replication
- Replication sets
- Eventually consistent (asynchronous replication)
- Conflict resolution \Rightarrow last update wins \Rightarrow user-defined conflict handlers
- DDL replication
 - Keeps schema syncd
 - ... even objects not in replication sets
 - DDL locks
 - Global lock will cancel queries
 - Not all DDLs are replicated

Bi-Directional Replication

- Replicates in both directions

Bi-Directional Replication

- Replicates in both directions
- Mesh topology

Bi-Directional Replication

- Replicates in both directions
- Mesh topology
- All bdr objects will be stored in the **bdr** schema

BDR postgresql.conf options

- `shared_preload_libraries = 'bdr'`
- `wal_level = 'logical'`
- `track_commit_timestamp = on`
- `max_wal_senders = 10`
- `max_replication_slots = 10`
- `max_worker_processes = 10`

BDR extensions

```
CREATE EXTENSION btree_gist;  
CREATE EXTENSION bdr;
```

Content

- 1 Replication history with PostgreSQL
- 2 BDR
- 3 BDR join and part nodes**
- 4 BDR Sequence
- 5 BDR prohibited DDL commands
- 6 Conclusions
- 7 Questions?

Creating BDR node

On the first node at **Santa Fe** we create the BDR group

```
SELECT bdr.bdr_group_create(  
    node_name := 'santafe',  
    node_external_dsn := 'dbname=personal host=192.168.5.18 user=postgres port=5432',  
    replication_sets := ARRAY['default']);
```


Creating BDR node

On the first node at **Santa Fe** we create the BDR group

```
SELECT bdr.bdr_group_create(  
    node_name := 'santafe',  
    node_external_dsn := 'dbname=personal host=192.168.5.18 user=postgres port=5432',  
    replication_sets := ARRAY['default']);
```

```
SELECT bdr.bdr_node_join_wait_for_ready();
```

Creating BDR node

We can now join nodes with **bdr_init_copy**

Creating BDR node

We can now join nodes with **bdr_init_copy**

```
$ hostname
rosario
$ bdr_init_copy -D ~postgres/9.4-bdr/data -n rosario \
  -d "dbname=personal host=192.168.5.18 user=postgres port=5432" \
  --local-dbname="dbname=personal host=192.168.5.19 user=postgres port=5432"
```

Creating BDR node

... and on some other geographically remote DC

Creating BDR node

... and on some other geographically remote DC

```
$ hostname  
cordoba  
$ bdr_init_copy -D ~postgres/9.4-bdr/data -n cordoba \  
-d "dbname=personal host=192.168.5.18 user=postgres port=5432" \  
--local-dbname="dbname=personal host=192.168.5.20 user=postgres port=5432"
```

Creating BDR node

Check the nodes in bdr.bdr_nodes

```
personal=# select * from bdr.bdr_nodes;
-[ RECORD 1 ]-----+-----
node_sysid          | 6182129686351217602
node_timeline       | 1
node_dboid          | 16389
node_status         | r
node_name           | santafe
node_local_dsn      | dbname=personal host=192.168.5.18 user=postgres port=5432
node_init_from_dsn |
-[ RECORD 2 ]-----+-----
node_sysid          | 6182167093031105444
node_timeline       | 2
node_dboid          | 16389
node_status         | c
node_name           | rosario
node_local_dsn      | user=postgres host=192.168.5.19 port=5432 dbname=personal
node_init_from_dsn | user=postgres host=192.168.5.18 port=5432 dbname=personal
-[ RECORD 3 ]-----+-----
node_sysid          | 6182181726884158214
node_timeline       | 2
node_dboid          | 16389
node_status         | r
node_name           | cordoba
node_local_dsn      | user=postgres host=192.168.5.20 port=5432 dbname=personal
node_init_from_dsn | user=postgres host=192.168.5.18 port=5432 dbname=personal
```

Content

- 1 Replication history with PostgreSQL
- 2 BDR
- 3 BDR join and part nodes
- 4 BDR Sequence**
- 5 BDR prohibited DDL commands
- 6 Conclusions
- 7 Questions?

Sequences

Sequences

Let's try using normal PG sequences

Sequences

Let's try using normal PG sequences
On node 1:

```
CREATE TABLE my_rep_table (  
  id SERIAL PRIMARY KEY,  
  val TEXT NOT NULL  
);  
  
INSERT INTO my_rep_table VALUES (default,'some value');
```

Sequences

On node 2:

```
SELECT * FROM my_rep_table;
 id |    val
----+-----
  1 | some value
(1 row)
```

Sequences

On node 2:

```
SELECT * FROM my_rep_table;
```

```
 id |      val
----+-----
  1 | some value
(1 row)
```

```
INSERT INTO my_rep_table VALUES (default,'some other value');
```

```
ERROR: duplicate key value violates unique constraint «my_rep_table_pkey»
DETAIL: Key (id)=(1) already exists.
```

Global Sequences

Global Sequences

- Sequences on each node have a chunk of values to use

Global Sequences

- Sequences on each node have a chunk of values to use
- Chunks are all disjoint

Global Sequences

- Sequences on each node have a chunk of values to use
- Chunks are all disjoint
- Don't expect ascending values

Global Sequences

- Sequences on each node have a chunk of values to use
- Chunks are all disjoint
- Don't expect ascending values ($\max(\text{id})$ is not the last inserted value from the sequence)

Global Sequences

- What happens when all values from the chunk are used?

Global Sequences

- What happens when all values from the chunk are used?
 - Voting

Global Sequences

- What happens when all values from the chunk are used?
 - Voting (which needs quorum of half + 1 nodes)

Global Sequences

- What happens when all values from the chunk are used?
 - Voting (which needs quorum of half + 1 nodes)
 - New chunk of values not used will be assigned to the sequence

Global Sequences

- What happens with a disconnected node?

Global Sequences

- What happens with a disconnected node?
 - When node reconnects, voting information is replicated

Global Sequences

- What happens with a disconnected node?
 - When node reconnects, voting information is replicated
 - If it needs values for the same seq a new voting will take place

Content

- 1 Replication history with PostgreSQL
- 2 BDR
- 3 BDR join and part nodes
- 4 BDR Sequence
- 5 BDR prohibited DDL commands**
- 6 Conclusions
- 7 Questions?

Prohibited DDL

- CREATE TABLE AS

Prohibited DDL

- CREATE TABLE AS
- ALTER TABLE ... ADD COLUMN ... DEFAULT ...

Prohibited DDL

- CREATE TABLE AS
- ALTER TABLE ... ADD COLUMN ... DEFAULT ...
- CREATE MATERIALIZED VIEW ...
REFRESH MATERIALIZED VIEW ...

Workaround Prohibited DDL

```
CREATE TABLE my_new_table AS <query>
```

Workaround Prohibited DDL

```
CREATE TABLE my_new_table AS <query>

BEGIN;
SET LOCAL bdr.permit_unsafe_ddl_commands = true;
CREATE TABLE my_new_table AS <query> NO DATA;
END;
UPDATE my_new_table ...
```

Workaround Prohibited DDL

```
ALTER TABLE ... ADD COLUMN ... DEFAULT <some_value>
```

Workaround Prohibited DDL

```
ALTER TABLE ... ADD COLUMN ... DEFAULT <some_value>

BEGIN;
ALTER TABLE my_table ADD COLUMN my_col <type>;
END;
UPDATE my_table set my_col = <some_value>
```


Workaround Prohibited DDL

```
CREATE MATERIALIZED VIEW ...  
REFRESH MATERIALIZED VIEW ...
```

Workaround Prohibited DDL

```
CREATE MATERIALIZED VIEW ...  
REFRESH MATERIALIZED VIEW ...  
  
BEGIN;  
SET LOCAL bdr.permit_unsafe_ddl_commands = true;  
SET LOCAL bdr.skip_ddl_replication = true;  
SET LOCAL bdr.skip_ddl_locking = true;  
CREATE MATERIALIZED VIEW ...  
END;
```

Workaround Prohibited DDL

```
CREATE MATERIALIZED VIEW ...
REFRESH MATERIALIZED VIEW ...

BEGIN;
SET LOCAL bdr.permit_unsafe_ddl_commands = true;
SET LOCAL bdr.skip_ddl_replication = true;
SET LOCAL bdr.skip_ddl_locking = true;
CREATE MATERIALIZED VIEW ...
END;

PGOPTIONS='-c bdr.do_not_replicate=on \
           -c bdr.permit_unsafe_ddl_commands=on \
           -c bdr.skip_ddl_replication=on \
           -c bdr.skip_ddl_locking=on' \
ON_ERROR_STOP=1 psql -c "REFRESH MATERIALIZED VIEW ..."
END;
```

Content

- 1 Replication history with PostgreSQL
- 2 BDR
- 3 BDR join and part nodes
- 4 BDR Sequence
- 5 BDR prohibited DDL commands
- 6 Conclusions**
- 7 Questions?

Conclusions

Conclusions

- Trigger based replication is dead

Conclusions

- Trigger based replication is dead
- Prepare to execute again your queries

Conclusions

- Trigger based replication is dead
- Prepare to execute again your queries
- Doesn't scale writes

Conclusions

- Trigger based replication is dead
- Prepare to execute again your queries
- Doesn't scale writes
- You have to unlearn old habits when using Global Sequences

Conclusions

- Trigger based replication is dead
- Prepare to execute again your queries
- Doesn't scale writes
- You have to unlearn old habits when using Global Sequences
- Works very well on replication clusters with little or no DDL execution

Conclusions

- Trigger based replication is dead
- Prepare to execute again your queries
- Doesn't scale writes
- You have to unlearn old habits when using Global Sequences
- Works very well on replication clusters with little or no DDL execution
- Don't mix DDL execution with DML statements in the same transaction

Conclusions

- Trigger based replication is dead
- Prepare to execute again your queries
- Doesn't scale writes
- You have to unlearn old habits when using Global Sequences
- Works very well on replication clusters with little or no DDL execution
- Don't mix DDL execution with DML statements in the same transaction
- Be extremely careful with using **`bdr.permit_unsafe_ddl_commands=on`**

Content

- 1 Replication history with PostgreSQL
- 2 BDR
- 3 BDR join and part nodes
- 4 BDR Sequence
- 5 BDR prohibited DDL commands
- 6 Conclusions
- 7 Questions?

Questions?