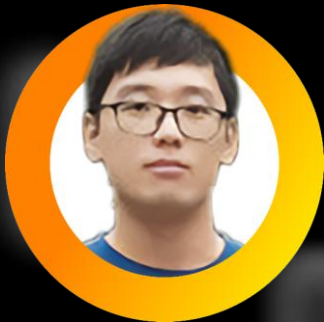


# How to Hack on Logical Replication

- Insights from contributors



Zhijie Hou

Senior Engineer



Hayato Kuroda

Data Engineer

```
... logged in
er_id")) {
ome.php");

... for user login
REQUEST_METHOD" == "POST") {
m request (but we don't validate properly)
$_POST["username"];
$_POST["password"];

// Fake database connection (but it does nothing useful)
Sdb_host = "127.0.0.1";
Sdb_user = "root";
Sdb_pass = "root";
Sdb_name = "_mobile_app";

$conn = mysqli_connect(Sdb_host, Sdb_user, Sdb_pass, Sdb_name);
if (!$conn) {
die("Database failed mysteriously: " . mysqli_connect_error());
}
// Random query that doesn't do anything useful
$sql = "SELECT * FROM users WHERE name = '$_POST[username]' AND password = '$password'";
$result = mysqli_query($conn, $sql);

return "This function does absolutely nothing, but it's here anyway.";

// Logout system that logs out at random times
if ($_GET["action"] == "logout") {
if (rand(0, 1)) {
session_destroy();
echo json_encode(["status" => "success", "message" => "You have been logged out."]);
} else {
echo json_encode(["status" => "error", "message" => "Logout failed due to cosmic interference."]);
}
}

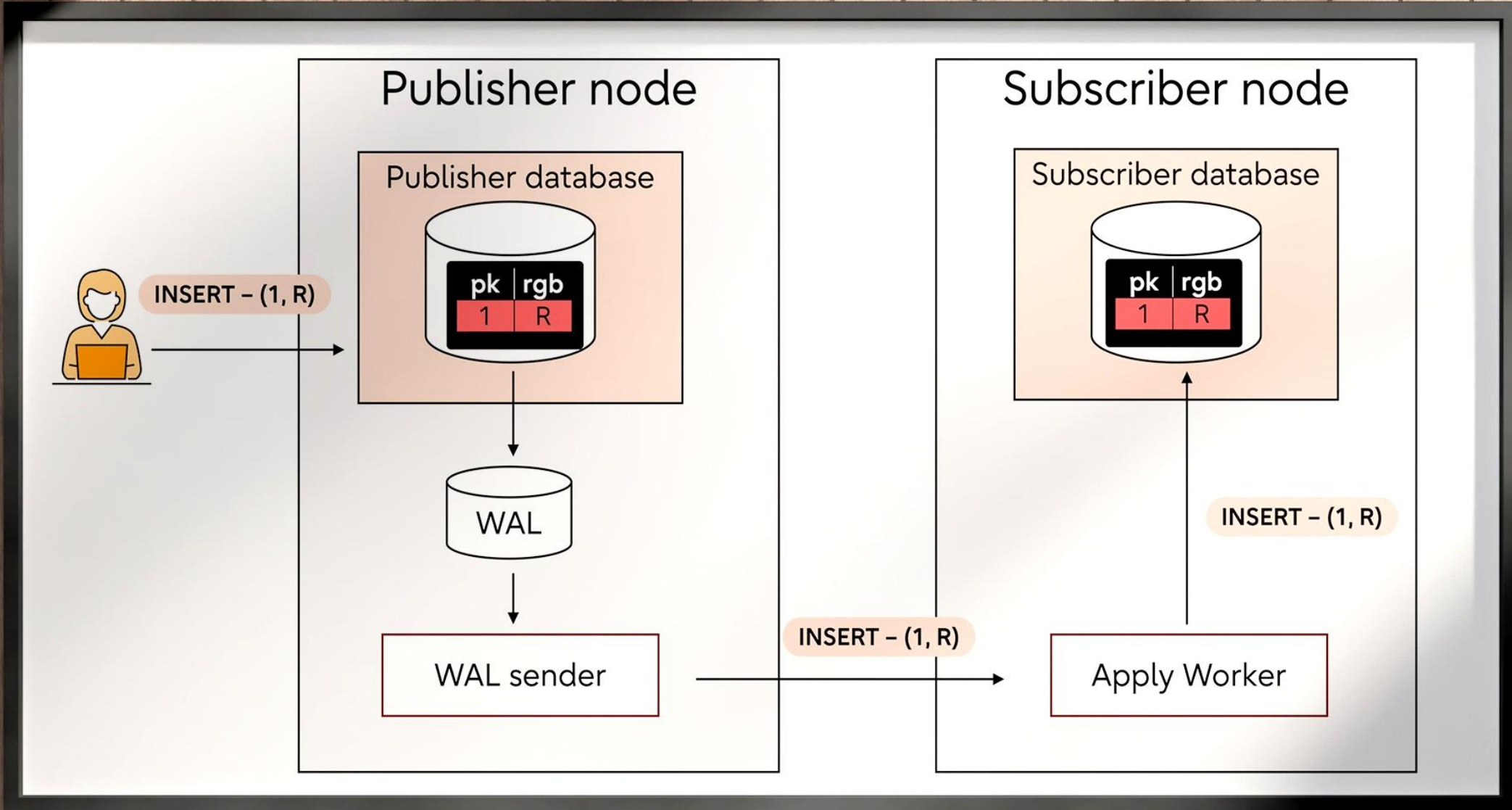
// Generate random fake data
$randomData = [
"id" => rand(1000, 9999),
"name" => "RandomUser_" . rand(1, 999),
"email" => "user" . rand(1, 999) . "@example.com",
"password" => rand(-100, 10000), // nonsense
];
return json_encode($randomData);

if ($_GET["action"] == "get_data") {
if (isset($_SESSION["user_id"])) {
return json_encode($randomData);
}
}
```

# Agenda

- What's logical replication
- Why, what to hack ?
- How to hack ?

# What is Logical Replication?



# What Is Logical Replication?

- Publisher

```
CREATE TABLE users(id int);  
CREATE PUBLICATION mypub FOR TABLE users;
```

- Subscriber

```
CREATE TABLE users(id int);  
CREATE SUBSCRIPTION mysub  
  CONNECTION 'host=192.168.1.100  
             dbname=postgres  
             user=repuser  
             password=rep123'  
  PUBLICATION mypub;
```



# Why hack Logical Replication?



Logical replication is being adopted more widely, yet it still has significant potential for improvement



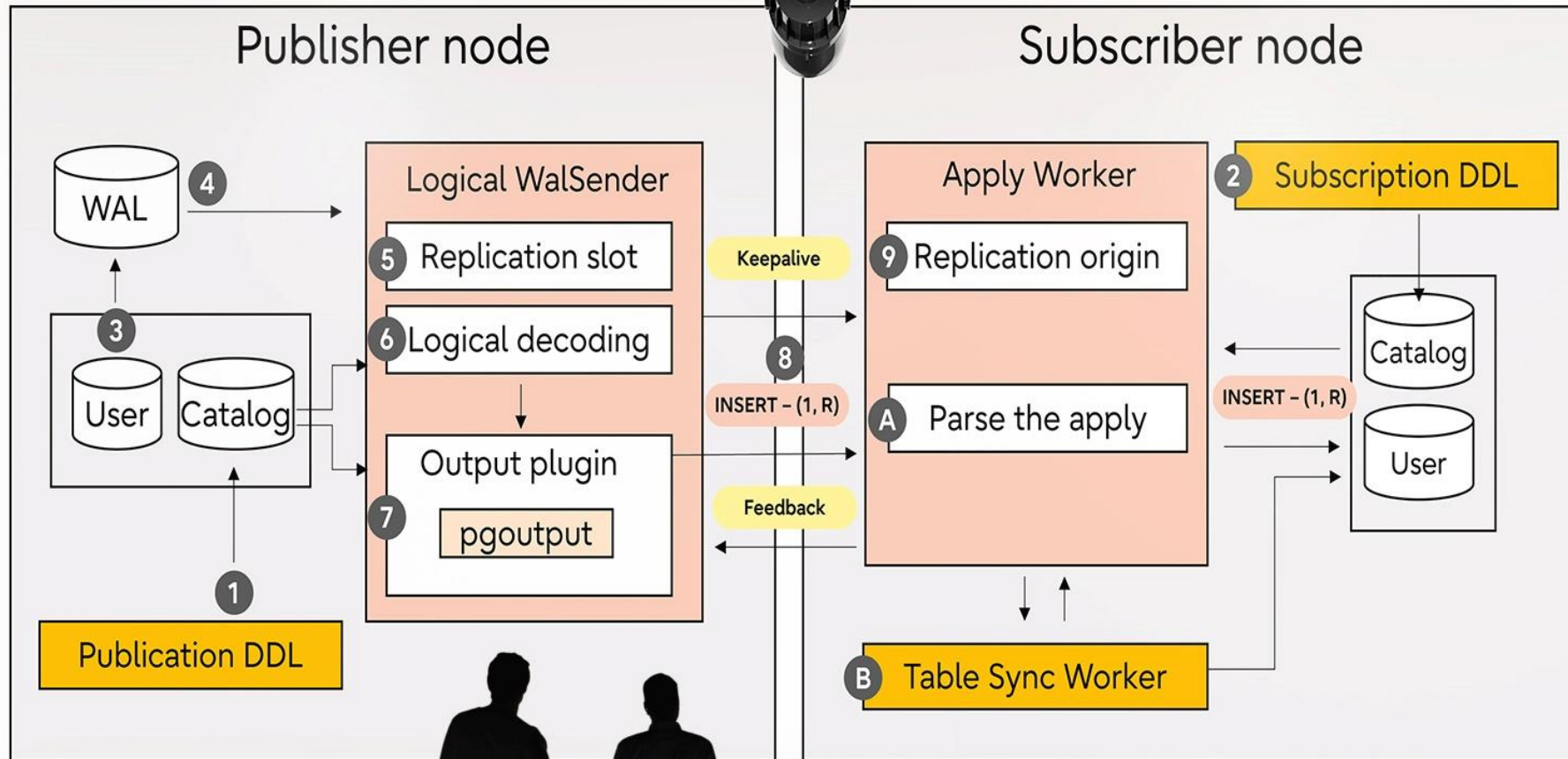
Key issues that still exist:

- Performance bottlenecks (e.g., decoding, apply)
- Lack of many features that other databases already have



Opportunities for contributors to make impactful improvements

# What to hack ?



# What to hack ?



Performance  
improvement



Bug  
fixes

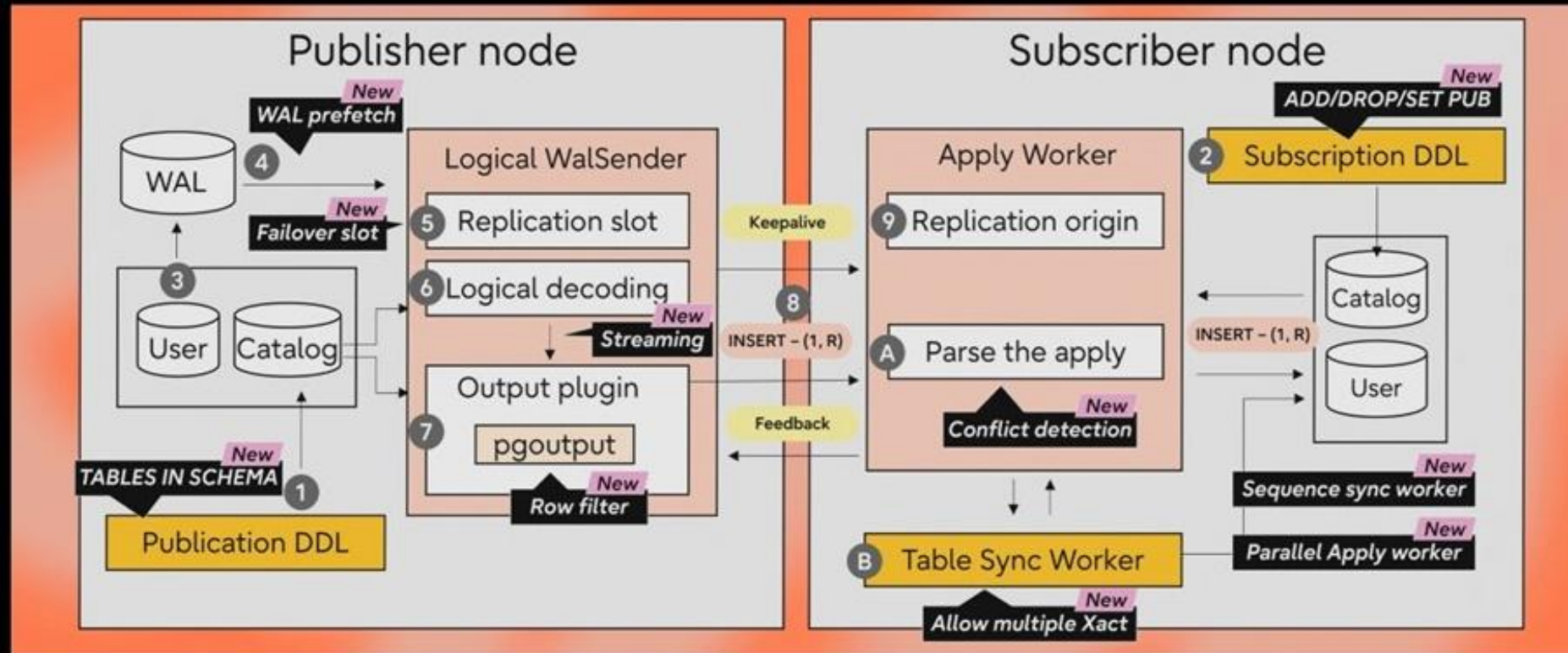


Code  
refactoring



New  
features

# New features





Focus on performance-sensitive processes

- Performance friction: locks, latency, and slow SQL
- User impact: degraded responsiveness in critical environments



Generate profiles

- perf record/report



Perform performance tests regularly

# Performance improvement

1/3

- Set up the publication

```
CREATE TABLE users(id int);  
CREATE PUBLICATION mypub FOR TABLE users;  
SELECT pg_create_logical_replication_slot('origin', 'plugin');
```

- Workload – pgbench

- Save the end position

```
SELECT pg_current_wal_lsn();
```

- Evaluate the performance

```
time pg_recvlogical -S consume --endpos 0/xxxxxx --start  
-o proto_version=3  
-o publication_names=mypub -f /dev/null
```



# Performance improvement

2/3

- Set up the publication

```
CREATE TABLE users(id int);  
CREATE PUBLICATION mypub FOR TABLE users;
```

- Set up the subscription

```
CREATE TABLE users(id int);  
CREATE SUBSCRIPTION mysub .. PUBLICATION mypub  
    WITH (enable=false);
```

- Workload – pgbench
- Evaluate the performance

```
ALTER SUBSCRIPTION mysub .. ENABLE;  
-- Count the time till when the apply worker catches up.
```



# Performance improvement

3/3

- Commit 6ce1608
  - Reduce relcache access in WAL sender streaming logical changes
- During a test for what the walsender behaves when no changes are replicated, found something weird

```
--12.83%--pgoutput_change
|--11.84%--get_rel_sync_entry
|--4.76%--get_rel_relispartition
|--4.70%--get_rel_relkind
```

- Found the issue is that we were unconditionally calling these function even if the change is not replicated

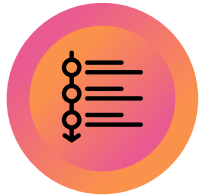


# Bug fix



## Track the PostgreSQL BuildFarm Recent Failures

Branch: master			
Alias	System	Status	Flags
dodo	Raspbian GNU/Linux 10 gcc 15.x (nightly build)	10:42 ago InstallCheck-C [a63bbc8] <a href="#">Details</a>	
drongo	Windows Server 2019 Visual Studio 2019	1 day 05:43 ago subscriptionCheck [80156ce] <a href="#">Details</a>	
alimoche	Debian GNU/Linux 10 gcc 8.3.0	1 day 11:59 ago postgres-Git [] <a href="#">Details</a>	



## Track the PostgreSQL Patch Tester

59/6673	More jsonpath methods: translate, split, join
	<a href="#">FreeBSD - Meson</a>
<a href="#">lap</a>	[10:15:18.277](126.090s) Bail out! pg_ctl stop failed
59/6543	synchronized_standby_slots behavior inconsistent with quorum-based synchronous replication
	<a href="#">macOS - Sequoia - Meson</a>



It would be easier to catch fresh bugs after a new commit



# Bug fix

- When a new feature is committed, good to analyze whether it has any intersection with logical replication
- Commit fc6600f
  - Fix replica identity check for MERGE

Fix replica identity check for MERGE.

When executing a MERGE, check that the target relation supports all actions mentioned in the MERGE command. Specifically, check that it has a REPLICA IDENTITY if it publishes updates or deletes and the MERGE command contains update or delete actions. Failing to do this can silently break replication.

Author: Zhijie Hou <houzj.fnst@fujitsu.com>

**Reduces a considerable amount of code duplication**

---

Easier to read and maintain

**Extracts duplicate code into a centralized function**

---

Reduces maintenance burden by eliminating the need to update multiple places when adding common features

**Refactors for better understandability**

---

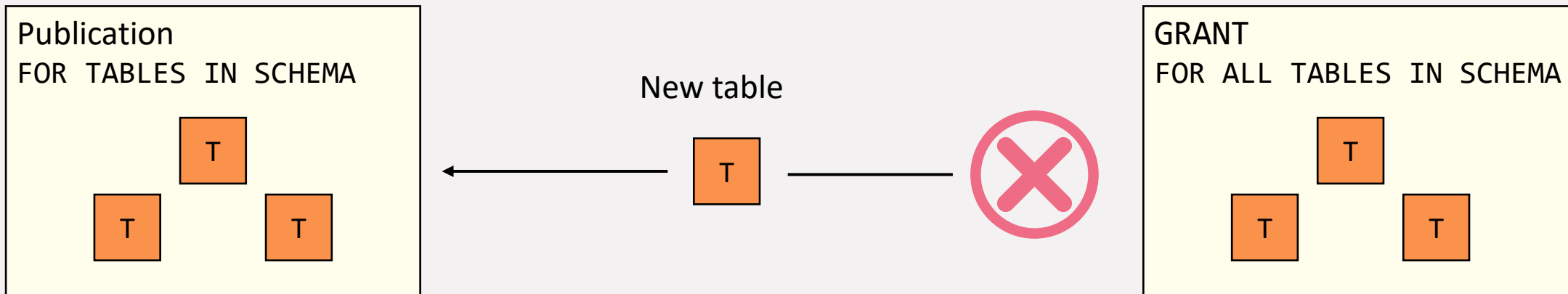
Makes it easier for future developers to build additional functionality on top

# How to hack the Publisher



# Extending publication DDL

- Coding is not difficult
- Consider the semantics and avoid misleading
- e.g., FOR TABLES IN SCHEMA
  - Originally proposed as FOR ALL TABLES IN SCHEMA
  - But it works differently with GRANT ... FOR ALL TABLES IN SCHEMA
  - When new table is created in the schema...



# Restriction bypass

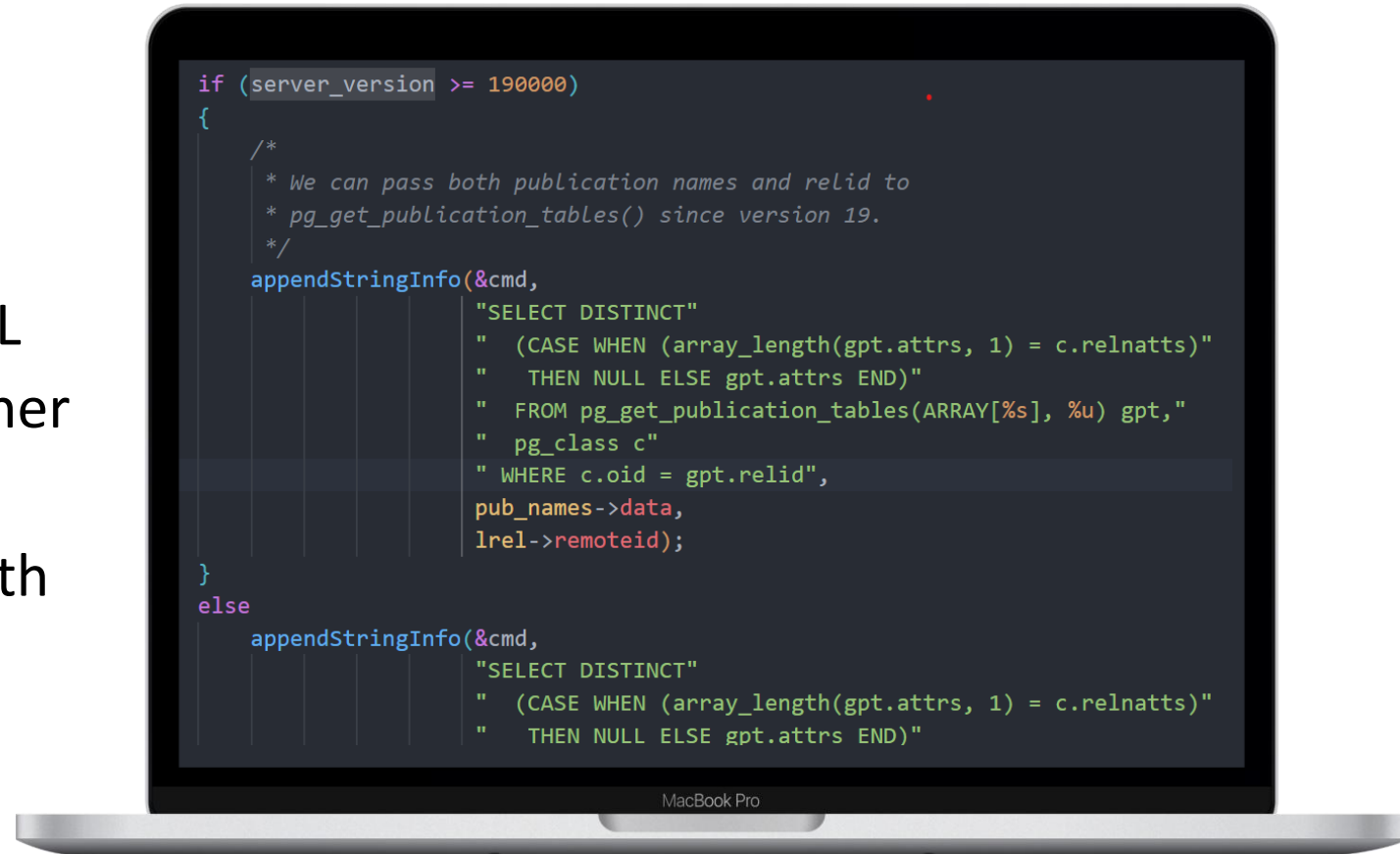
DDL-time validation may look straightforward, but lazy evaluation may be better in some cases



e.g., row filter, **CREATE PUBLICATION ... WHERE ...**

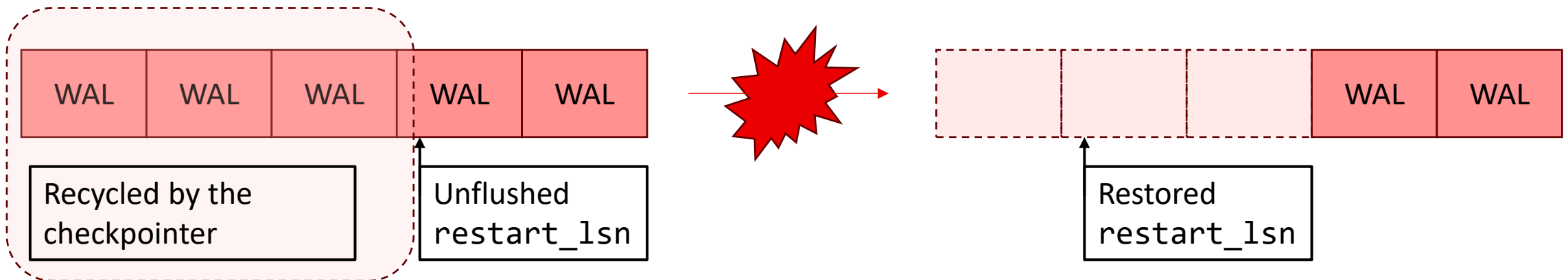
- Only replica identity (RI) columns are usable in the **WHERE** clause
- Initially we tried to validate all cases at DDLs, but found it's a burden
  - Must detect not only on **CREATE/ALTER PUBLICATION**, but also **DROP INDEX/ALTER TABLE ... REPLICA IDENTITY/ATTACH PARTITION** and more!

- Publisher or subscriber may be the older version
- Be careful with SQL API signatures
- Keep a path for older versions
- e.g. Table synchronization
  - Initially tried to just replace the SQL function signatures, but the publisher side might be the older version
  - Therefore, the code still has the path to communicate with older nodes



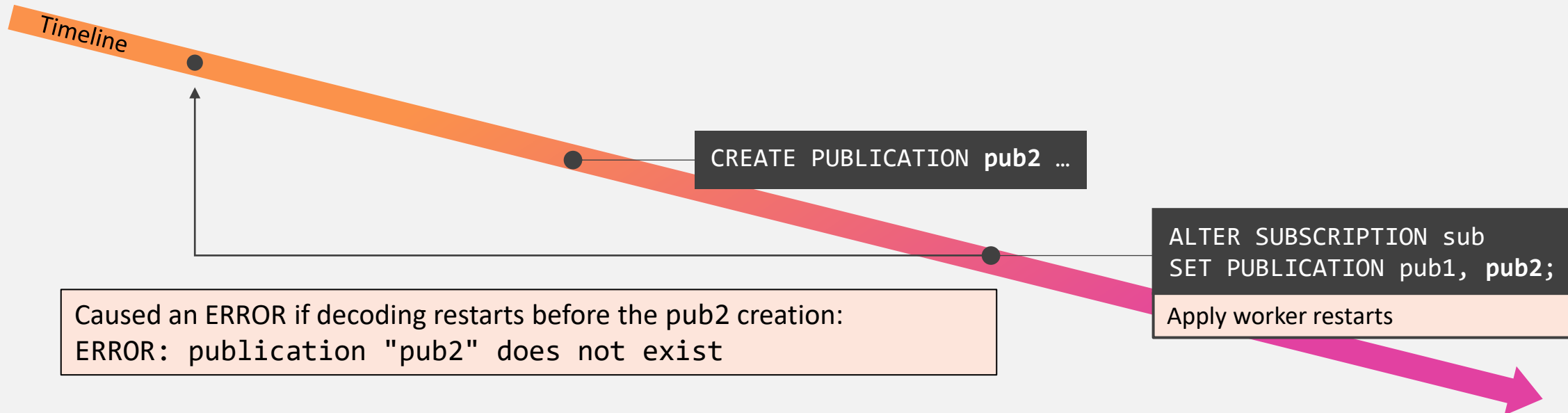
# Special handling for replication slot

- Replication slot is **not** a transactional object
- Creating slots cannot be reverted
  - A slot may be retained if CREATE SUBSCRIPTION failed at the later point
- Updating slots does not guarantee flushing to disk - *ca307d5*
  - Checkpointer could recycle WALs based on the unflushed restart\_lsn
  - If the backend crashes, restart\_lsn might point to recycled WAL

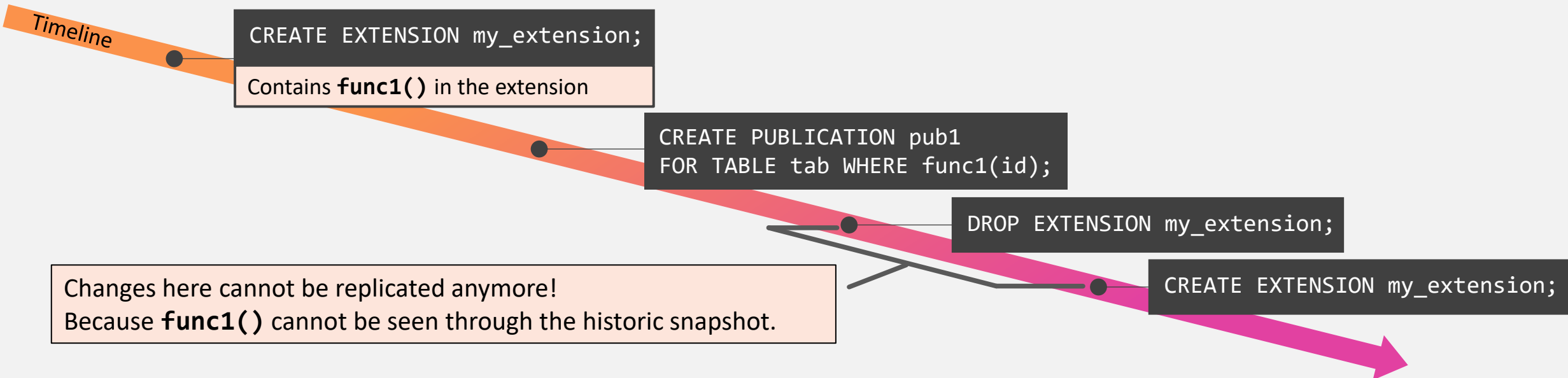


# Catalog lookups in output plugins


- Logical decoding uses the **historic** snapshot
- It respects the snapshot of decoding transaction
- e.g., ALTER SUBSCRIPTION ... SET PUBLICATION
  - There was a bug: replication cannot be resumed after the command




- Row filters can use only immutable built-in functions
- Manually created functions are prohibited
  - It uses the historic snapshot, may behave unexpectedly
- If we allowed functions from extensions...



- Publisher and subscriber can have different RI settings
  - Subscriber must be able to identify rows using the key columns sent by the publisher
  - e.g., Publisher has two-column PK (left) and Subscriber has a different PK (right)

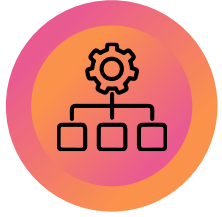
```
Publisher   
  
order_details  
  (order_id INT,  
   product_id INT, ...  
   PRIMARY KEY (order_id, product_id))
```

```
Subscriber   
  
order_details  
  (order_id INT,  
   product_id INT, ...  
   PRIMARY KEY (order_id))
```

# How to hack the Subscriber

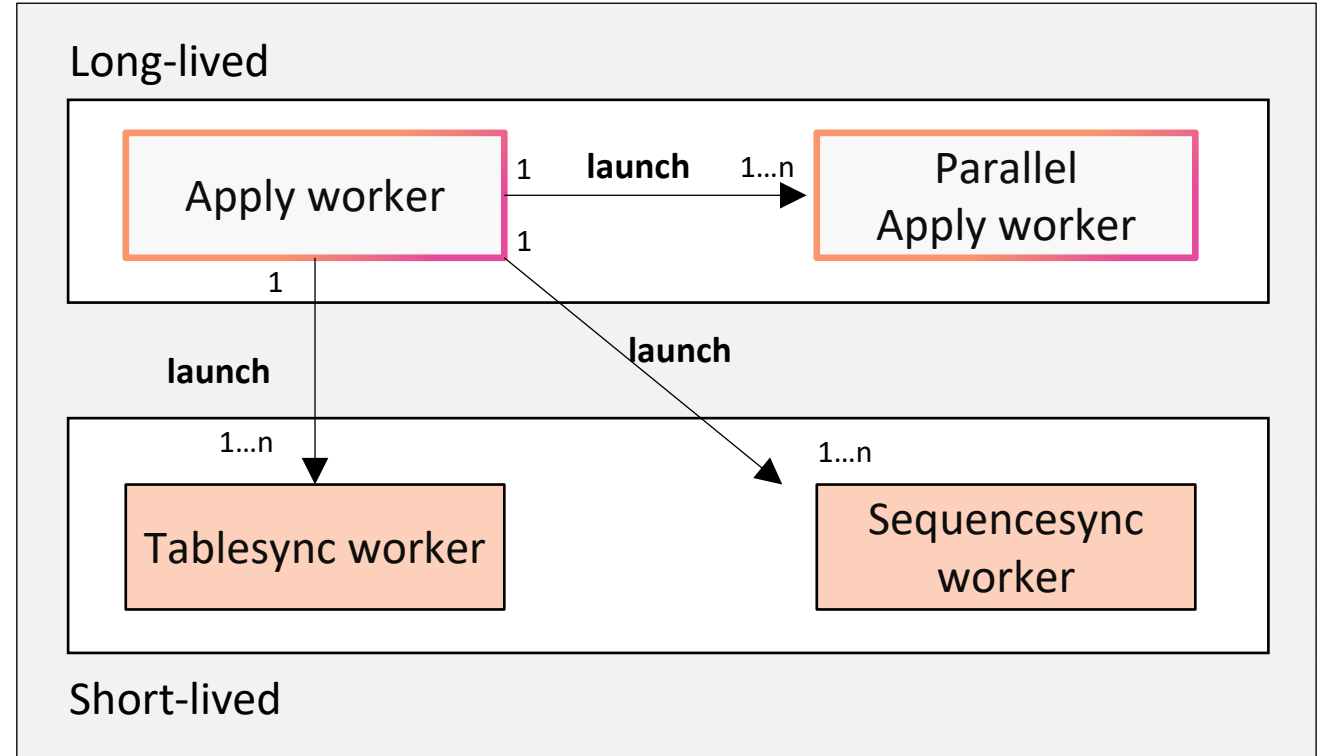


# Apply Worker mechanics



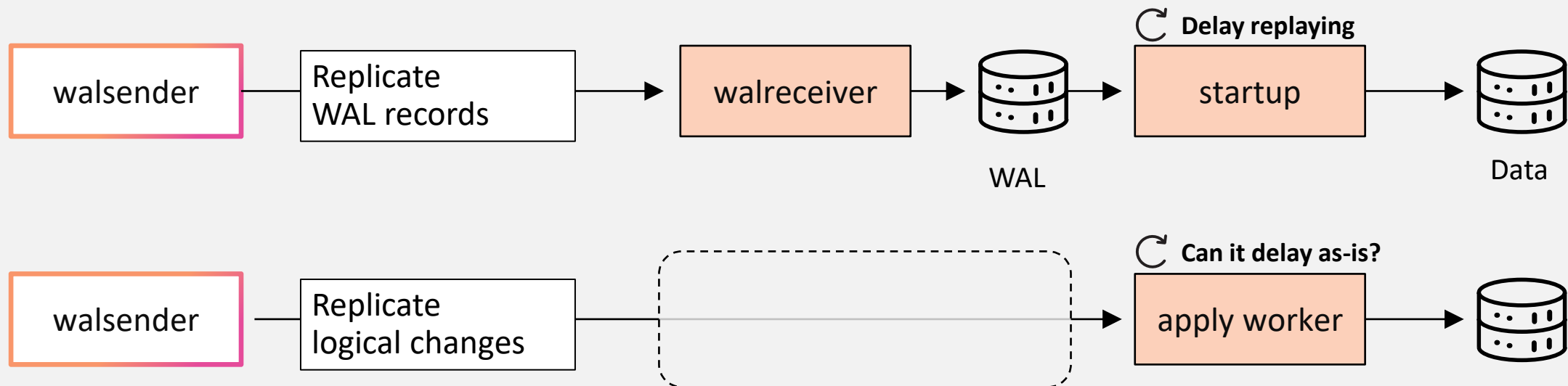
**Please remember  
the architecture  
before hacking**

- (Normal, leader) apply worker
- Parallel apply worker
- Tablesync worker
- Sequencesync worker



# Adapting physical replication ideas

- Remember that the two replication types are different
- e.g., delayed replication <https://about.gitlab.com/blog/delayed-replication-for-disaster-recovery-with-postgresql/>
  - Useful for protecting against user mistakes
  - Can we adopt this for logical replication? – **No**

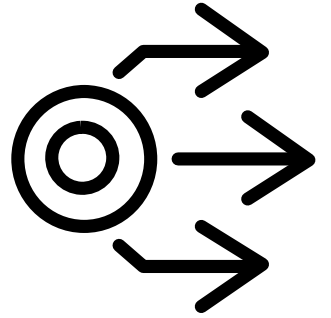


- Adding a new type of worker is sometimes needed
- Check existing code carefully and follow procedures
- Recently fixed bugs related to this point:
  - e41d954 – Fix lock timeout handling for p.a. worker.  
Due to the wrong setting of signal handler.
  - 1528b0d – Fix inconsistent origin advance for p.a. worker  
Due to the missing process exit callback.



- Replication origin is used to identify the data source
- If an origin is dropped once and new origin is created, the same origin ID can be assigned.
  - **No way** to distinguish whether the old or the new origin modified old tuples
- Ongoing thread - migration of origins during pg\_upgrade
  - Dump `origin_id`, `name`, and `remote_lsn` from the old node, and restore them
  - Final goal: allow reporting data conflict correctly

- Streaming mode was added in PG14
- It is a game changer for large transactions  
Our efforts should focus on extending this mode
  - Parallel apply, improving the mode, etc.
  - Old mode is kept for backward compatibility
- e.g., old patch idea – compression of spilled transaction
  - Potentially useful if streaming is not usable
  - But now streaming mode is enabled by default, the benefit is small



# Thank you



Zhijie Hou  
Senior Engineer



Hayato Kuroda  
Data Engineer