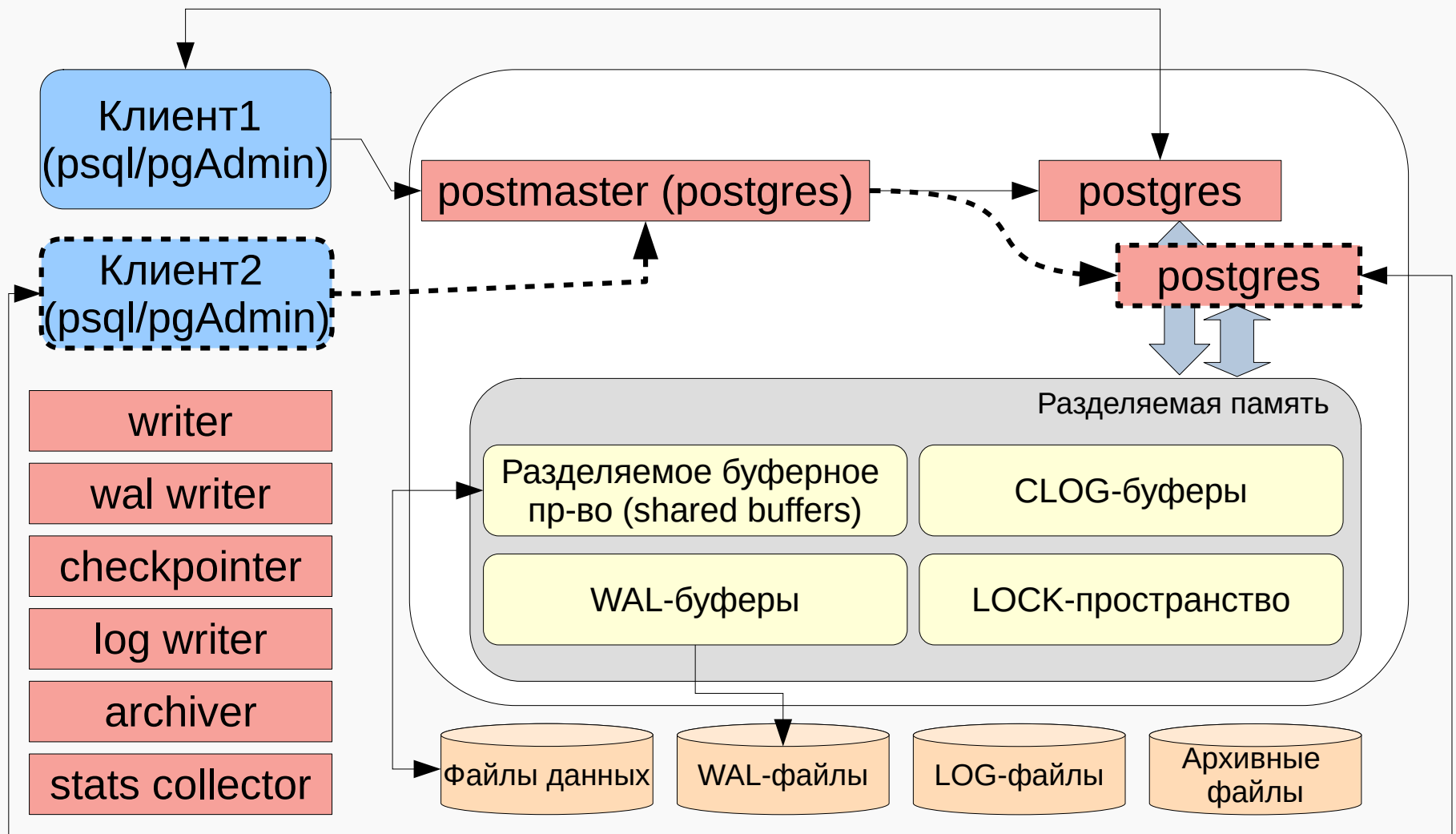


# Лекция 8

Лектор:  
Николаев В.В.

# 1. WAL в PostgreSQL

# Архитектура PostgreSQL



# WAL buffers

WAL — Write Ahead Log:

- Хранит информацию об изменениях данных в БД — WAL-записи (XLOG).
- Эта информация используется для воссоздания актуального состояния данных в случае восстановления базы данных (например, после сбоя).
- Настраивается через параметр `wal_buffers` — размер wal-буфера:
  - по умолчанию  $1/32$  от разделяемого буфера;

# WAL-записи

- WAL-записи фиксируют различные изменения состояния данных в БД:
  - INSERT, UPDATE, DELETE, COMMIT, ...;
- WAL-записи записываются в:
  - Write Ahead Log;
  - При операции COMMIT — происходит синхронизация с WAL (логом транзакций);
- WAL — общий для всего кластера.

# PostgreSQL: WAL-записи

- Поддерживается REDO-лог.
- Запись об изменении попадает в журнал (на диск) раньше, чем данные на диск и clog.
- PGDATA/pg\_wal:
  - журнал разбит на сегменты;
  - wal-segsize — размер сегмента;
  - wal-segsize можно изменить только при инициализации кластера PostgreSQL;

# Особенности работы журнала

- В журнал не попадают:
  - действия над временными таблицами;
  - действия над UNLOGGED-таблицами;
  - **до** версии 10: хэш-индексы.
- UNLOGGED TABLE: изменения для данных в такой таблице не хранятся в WAL:
  - CREATE **UNLOGGED** TABLE Students ...

UNLOGGED можно использовать, **только** когда данные в таблице **не важны** (не страшны потери), и при этом требуется минимизировать число обменов IO

# PostgreSQL: WAL-записи

Основные составляющие записи в журнале:

- Заголовок:
  - размер всей записи;
  - номер транзакции;
  - размер данных;
  - ссылка на предыдущую запись;
  - контрольная сумма;
- Данные об изменении — могут быть по-разному представлены.

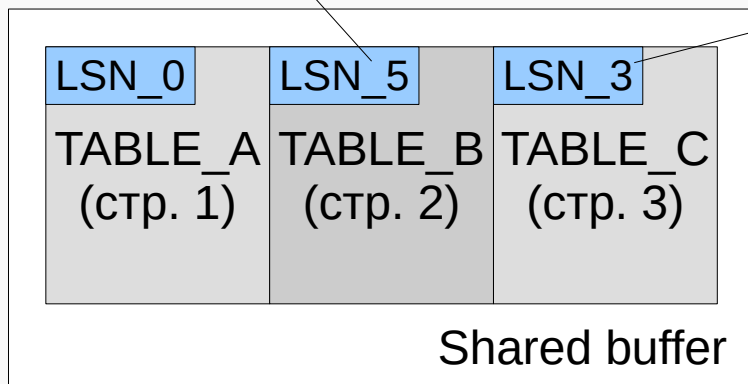


# WAL-записи

- У каждой WAL-записи есть **LSN (Log Sequence Number)** — уникальный идентификатор WAL-записи.
- Используется для нахождения позиции записи в WAL:
  - LSN — смещение до записи от начала журнала;
- LSN последней WAL-записи, относящейся к некоторой странице хранится в заголовке страницы.
- Для определения того, актуальна или нет некоторая страница (например, на момент восстановления) её LSN сравнивается с LSN в журнале.

# WAL-записи

У страниц есть заголовок.



В заголовке параметр `pd_lsn`:  
**последний** LSN с модификацией в этой странице

# Восстановление

- WAL (**REDO-лог**) — журнал последовательно читается и перевыполняются операции:
  - для определения записей, которые нужно заново «выполнить» используется LSN.
- Чтобы сократить число записей для просмотра: используются **контрольные точки**:
  - процесс **checkpoint**;
  - начало последней контрольной точки — REDO point.

# wal\_level

Параметр позволяет задать уровень детализации операций в WAL:

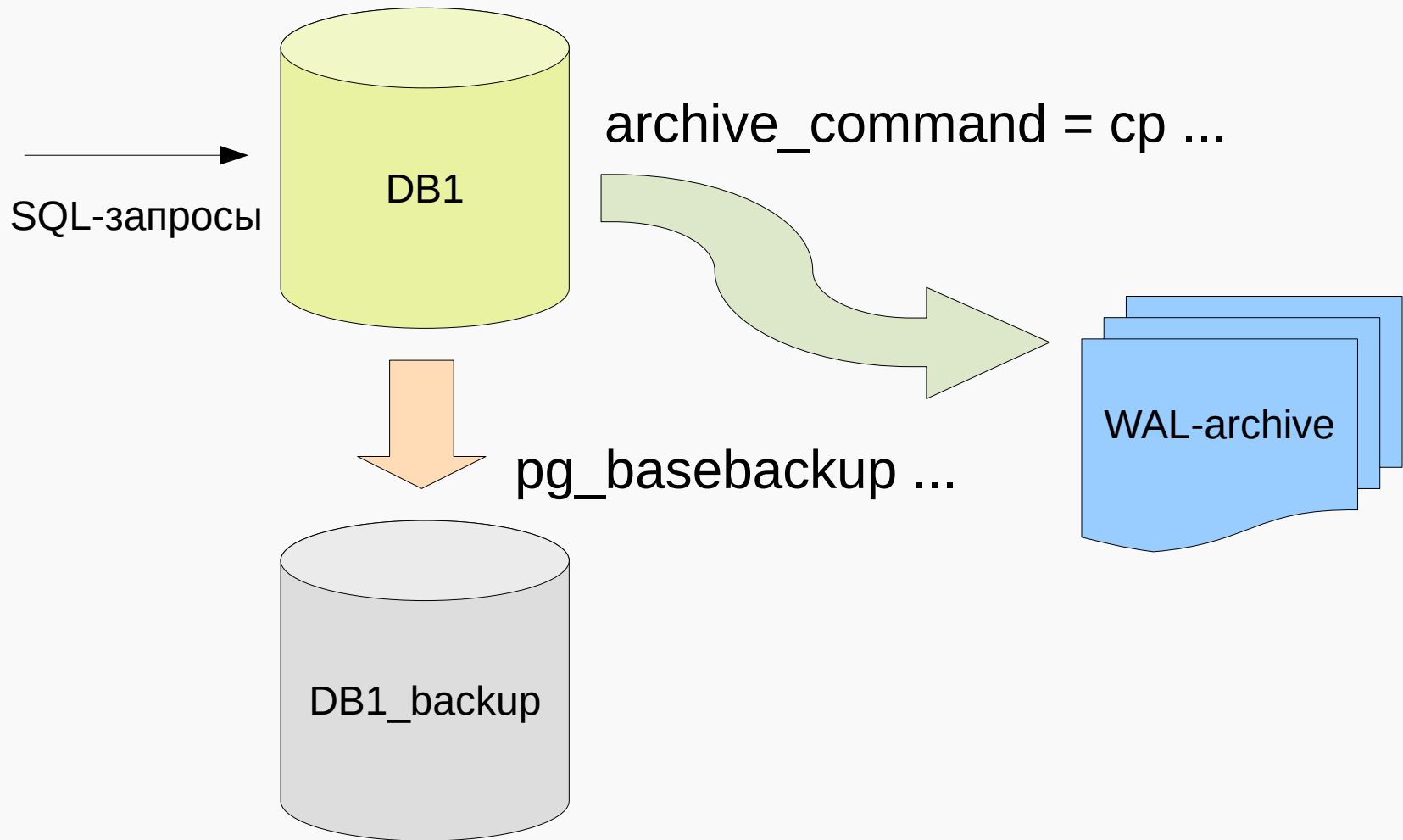
- **minimal** — только для восстановления после сбоя или immediate shutdown;
- **replica** — minimal + поддержка архивирования WAL и физической репликации;
- **logical** — replica + поддержка логической репликации (логического декодирования WAL);

# Непрерывное архивирование

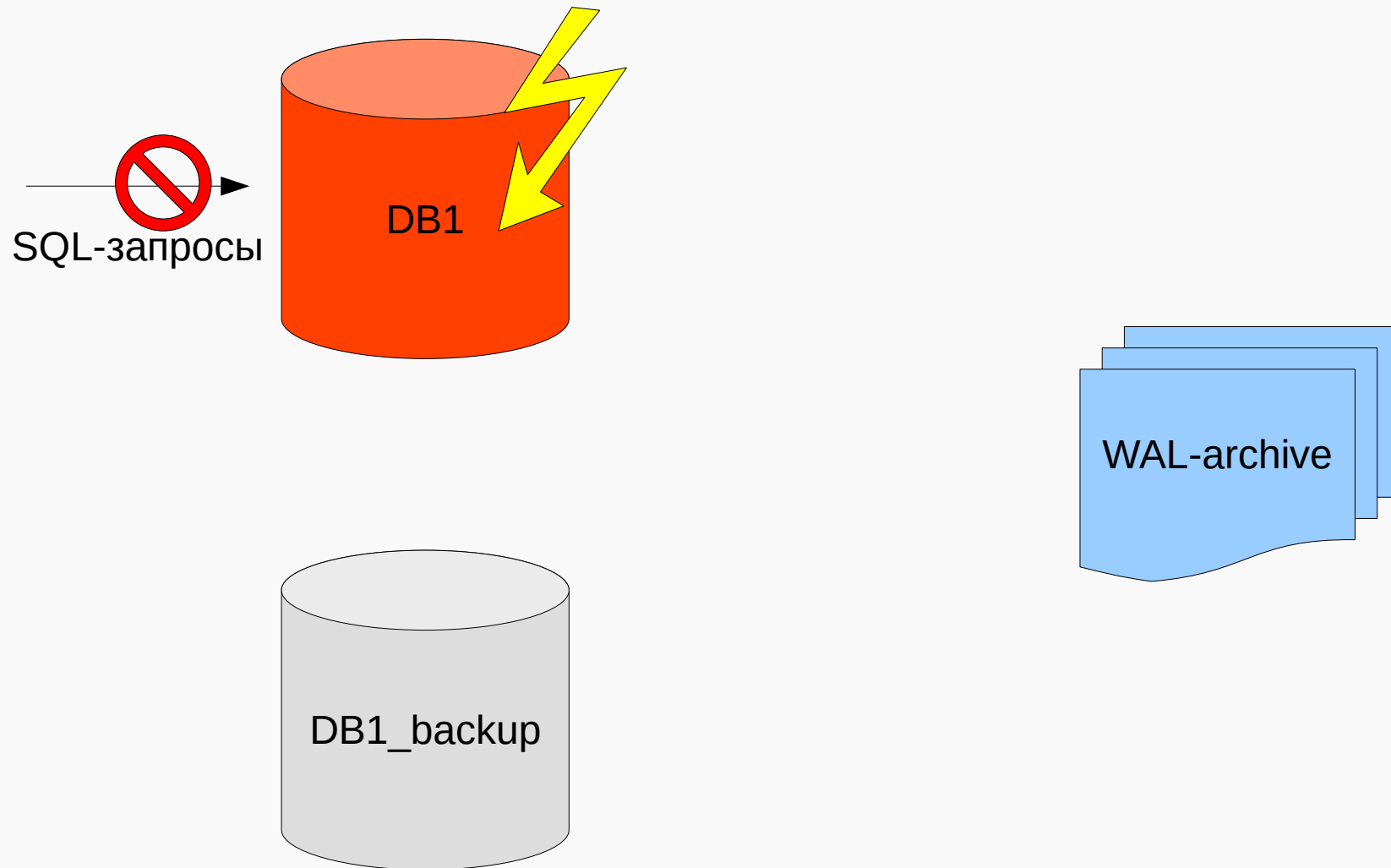
- Базируется на существовании WAL (PGDATA/pg\_wal):
  - состояние БД может быть восстановлено путем повтора зафиксированных в логе операций;
  - конфигурационные файлы не восстанавливаются через WAL (postgresql.conf, ...);

**Резервная копия = полная копия + WAL-файлы**

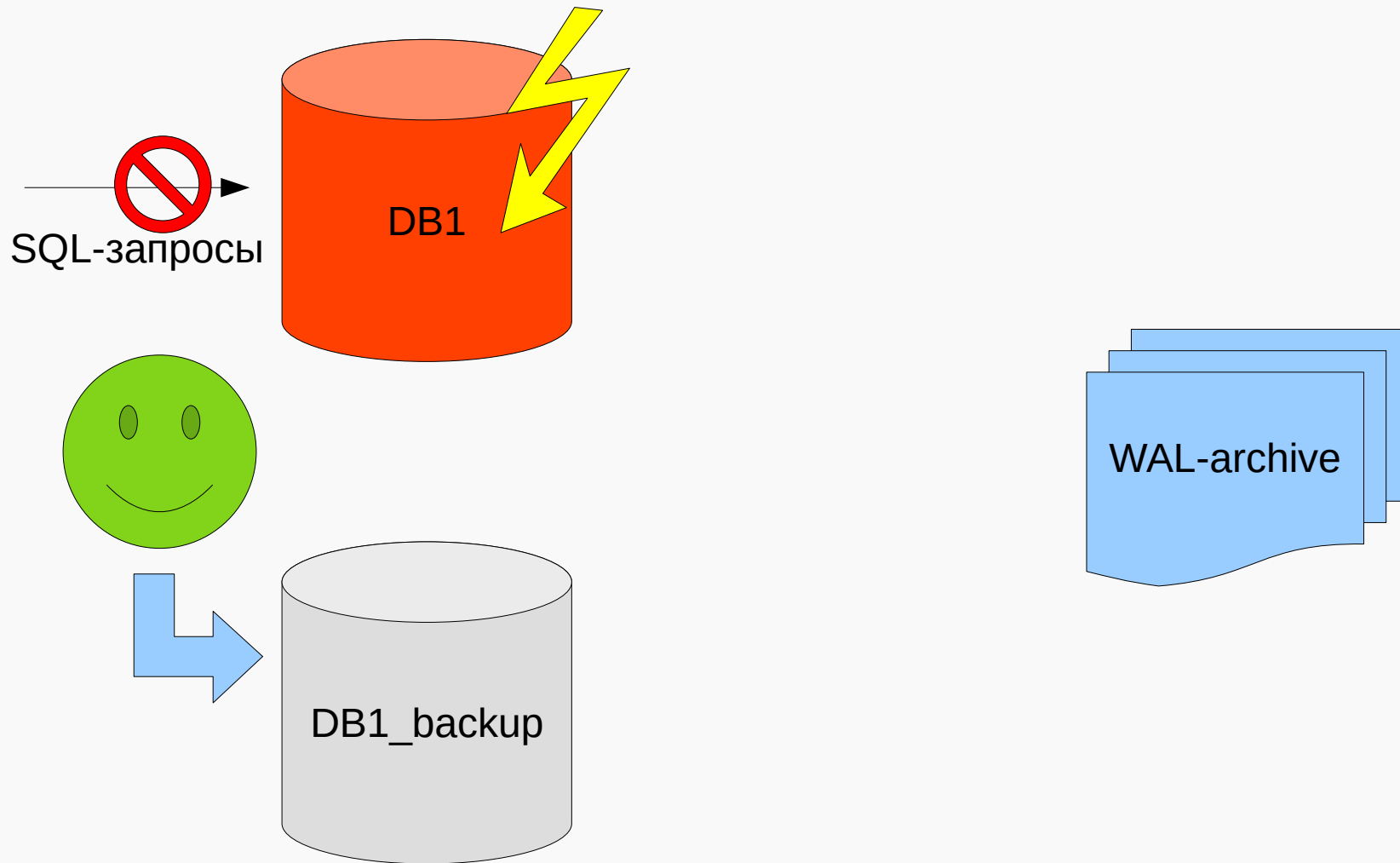
# 1. Нормальный процесс работы при непрерывном архивировании



## 2. Сбой

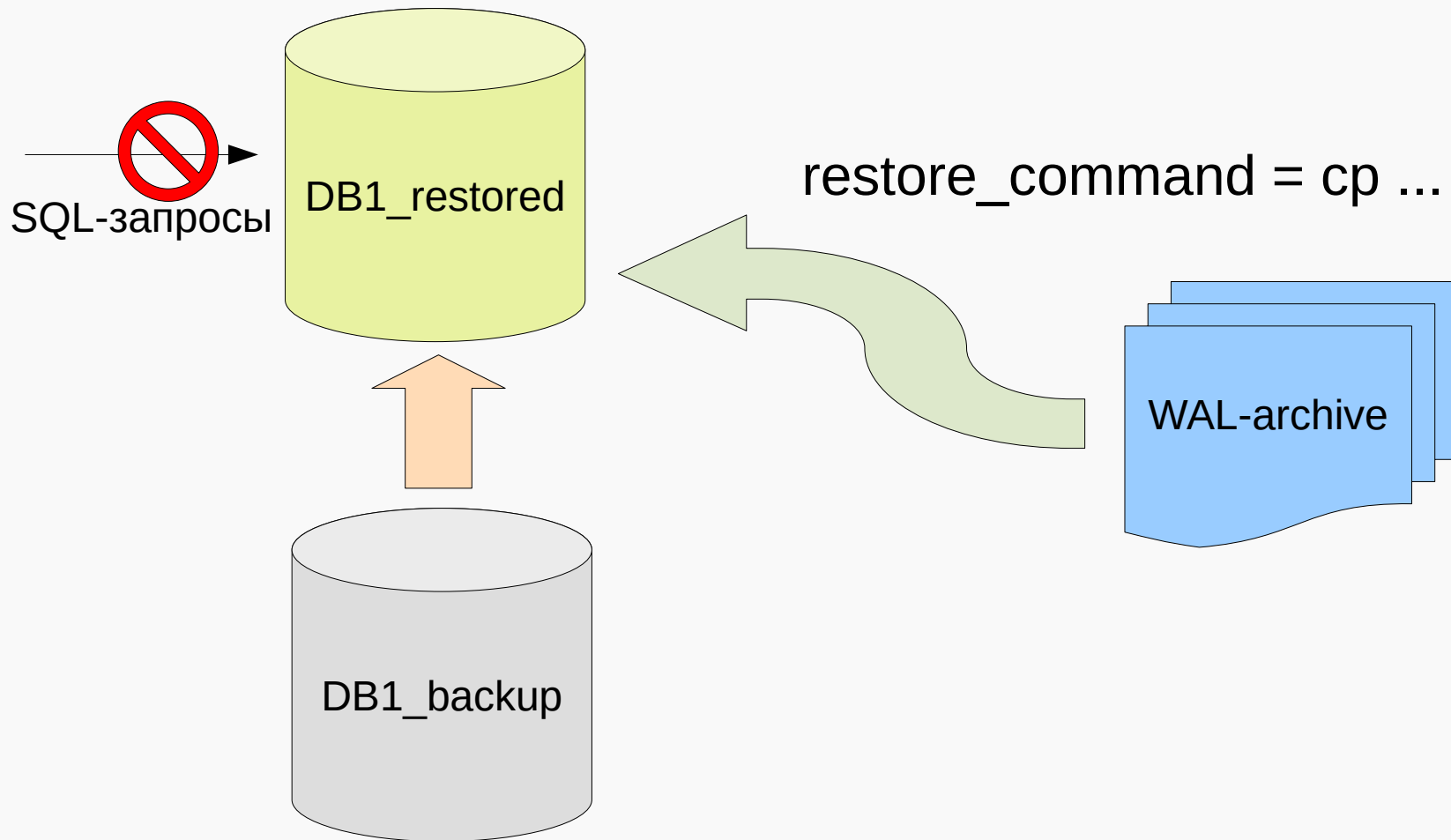


# 3. Обнаружение проблемы

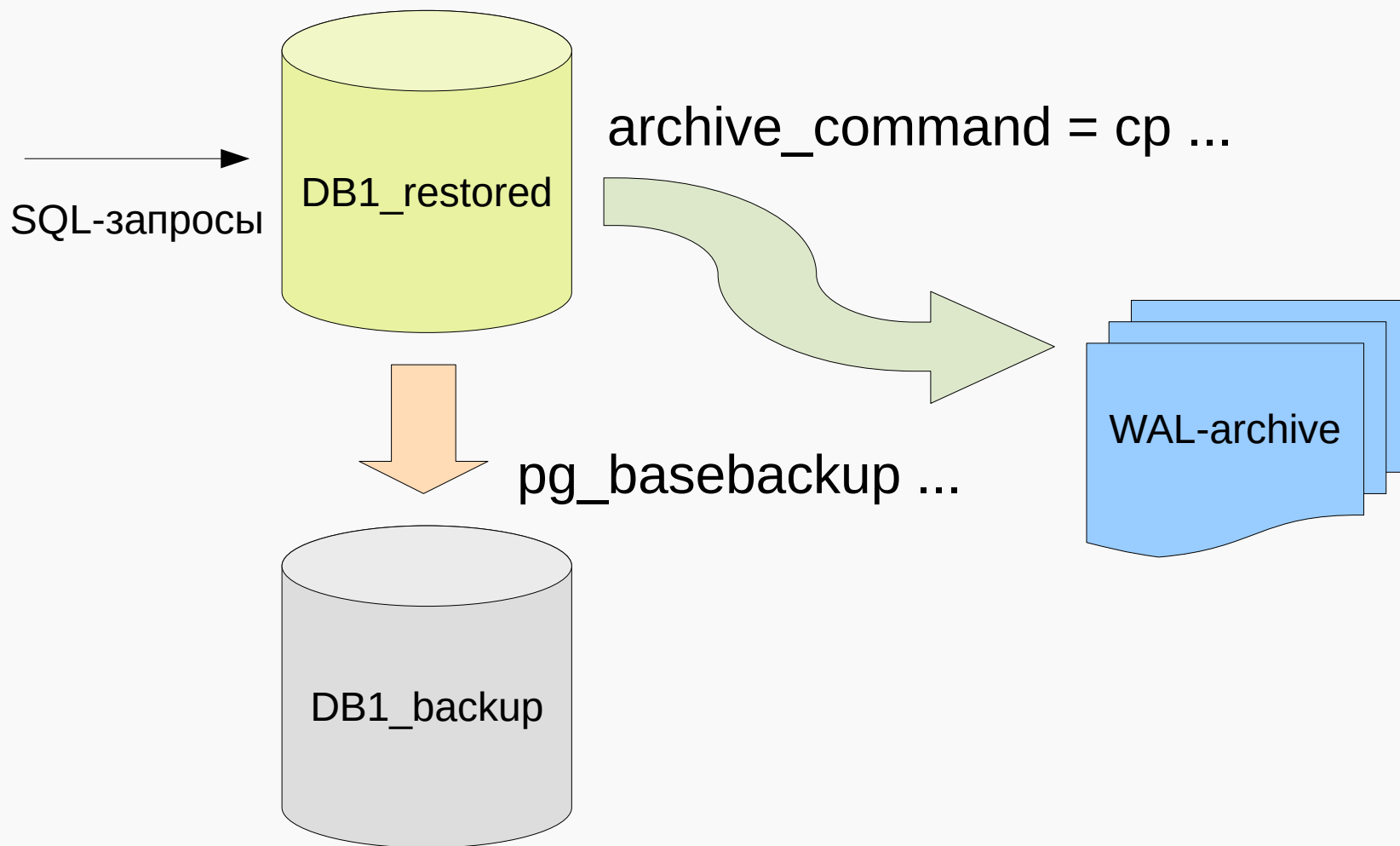




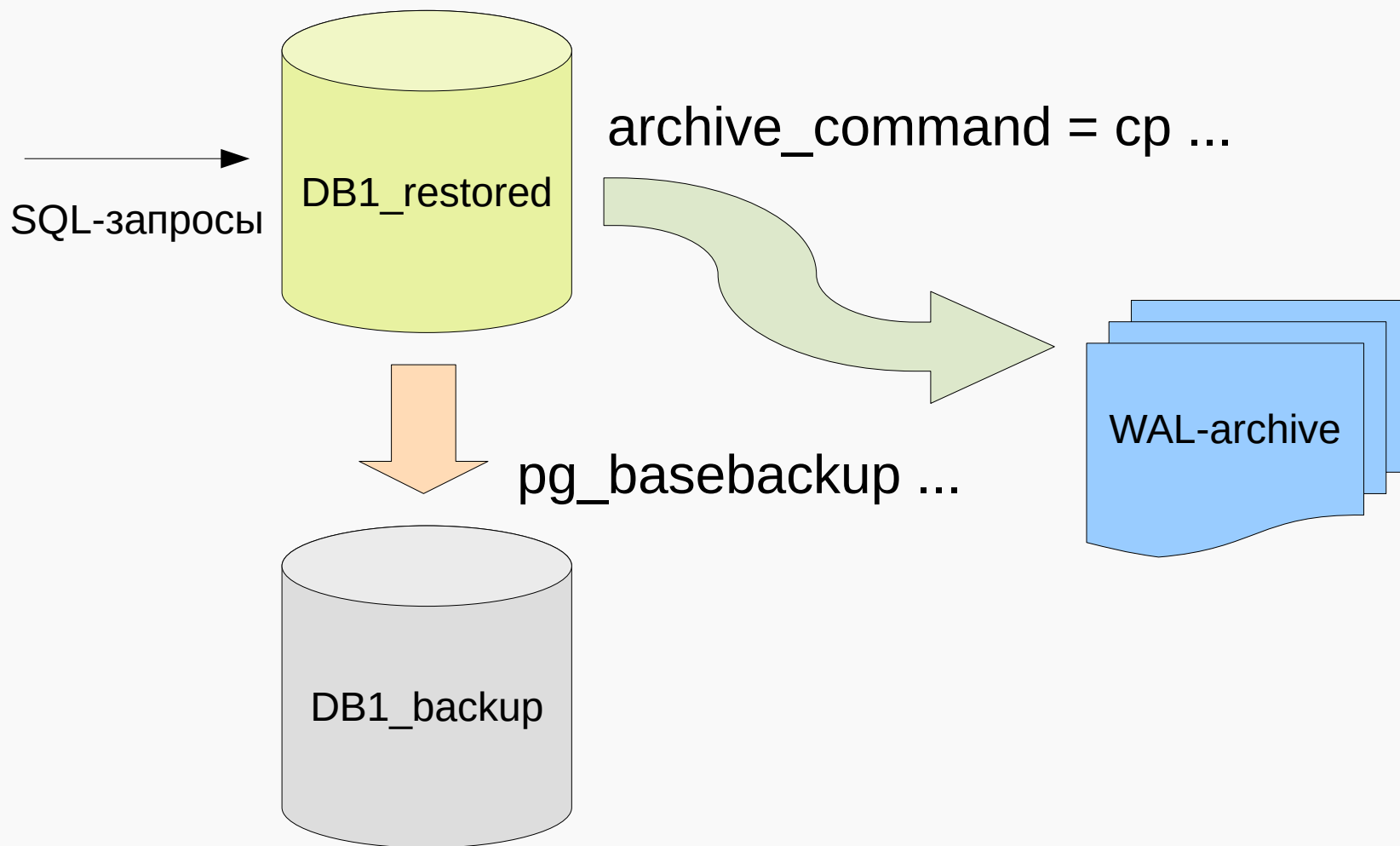
## 4. Восстановление



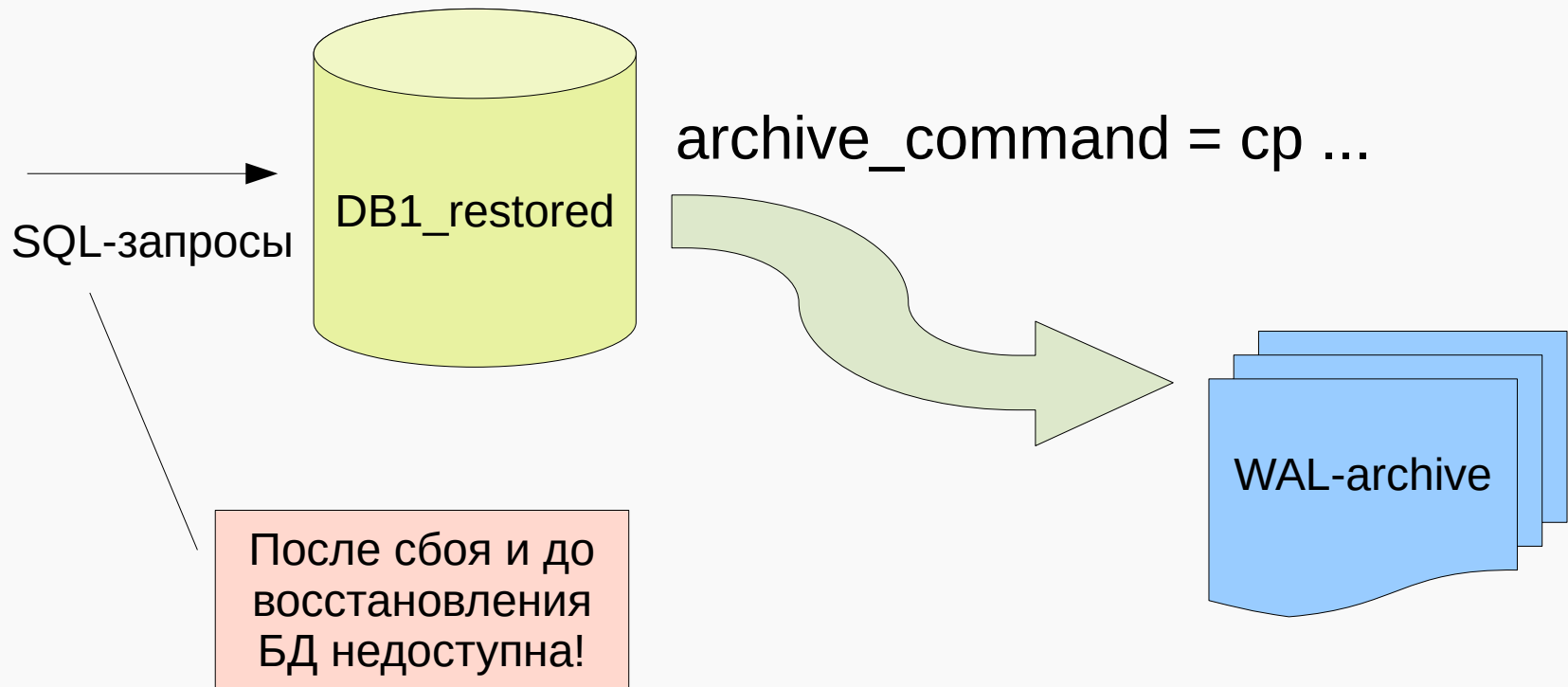
## 5. Состояние восстановлено



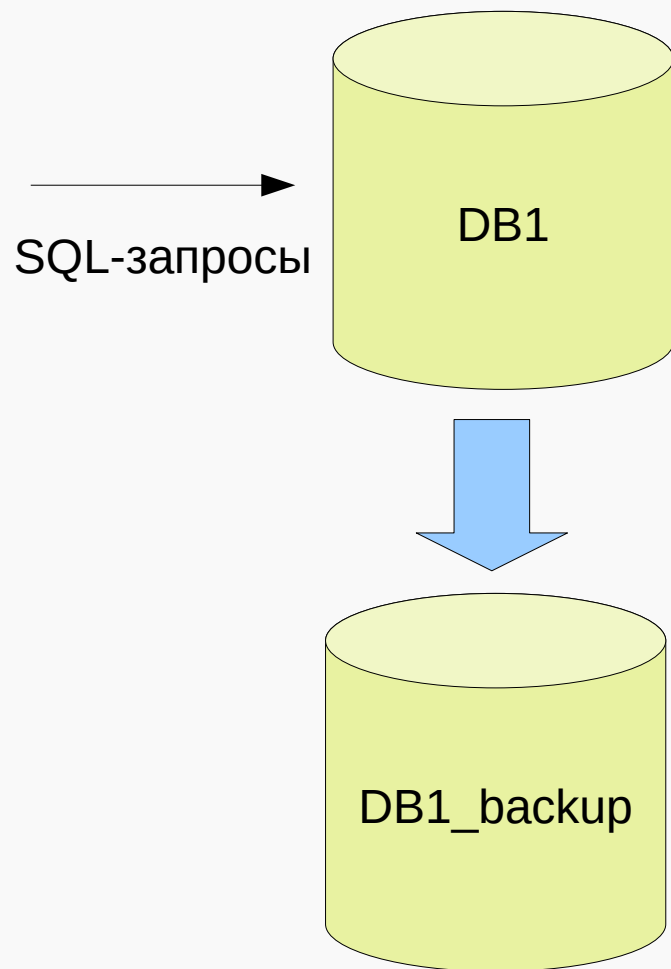
# В чем проблема?



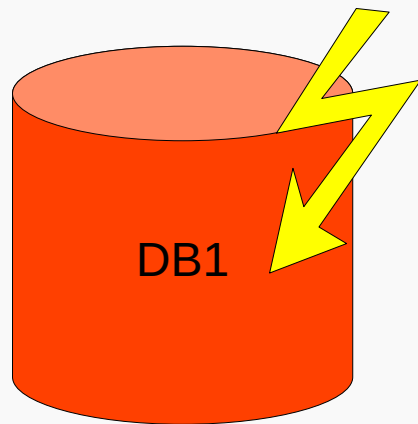
# В чем проблема?



# Возможное решение



# Возможное решение:



→  
SQL-запросы



## 2. Репликация данных в PostgreSQL

# ЧТО МЫ ХОТИМ?

**Высокий уровень доступности** — (high availability) способность системы работать без сбоев (непрерывно) в течение определенного периода времени.

Как обеспечить высокий уровень доступности:

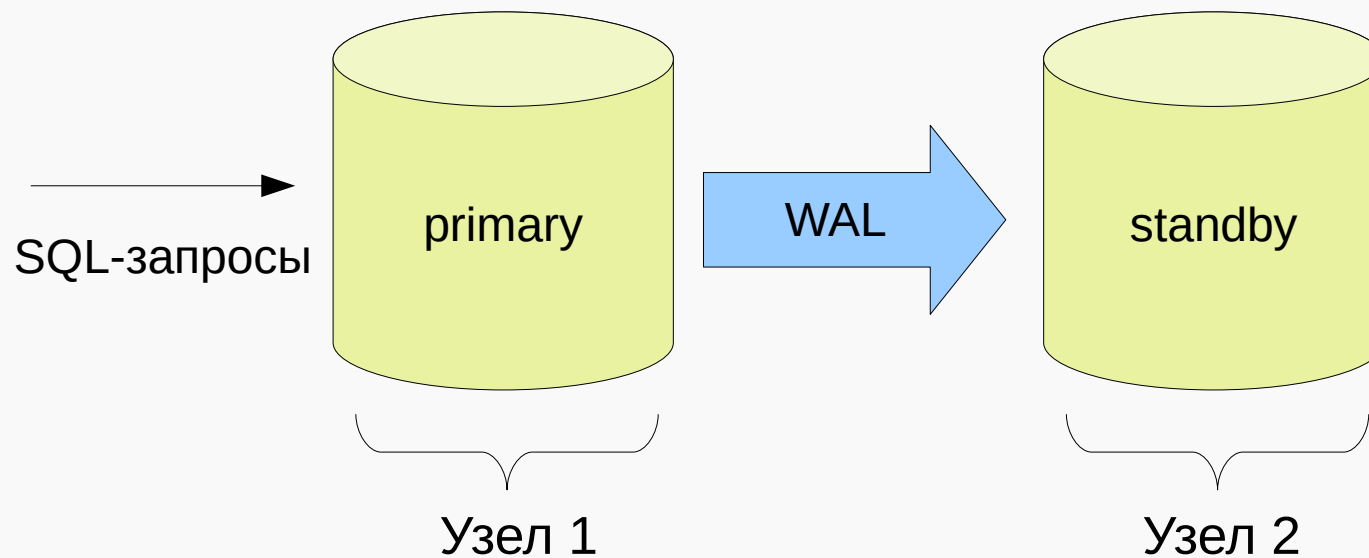
- обеспечить избыточность важных частей системы;
  - репликация данных.
- Failover, switchover.



# Репликация

- Репликация данных — средство обеспечения избыточности. Заключается в сохранении и поддержании нескольких копий данных.
- Репликация в СУБД:
  - **Узлы** — отдельные серверы БД, обеспечивающие работу. Совокупность узлов — **кластер**.
  - У каждого узла — своя БД (оригинал или копия, реплика). Копии данных синхронизируются в реальном времени.
  - Один из серверов БД — главный (**master, primary**). Режим — чтение/запись. Если несколько master-узлов — multimaster-репликация.
  - Остальные узлы — зависимые (**slave, standby**). Если режим — только чтение — **hot standby**.

# Репликация



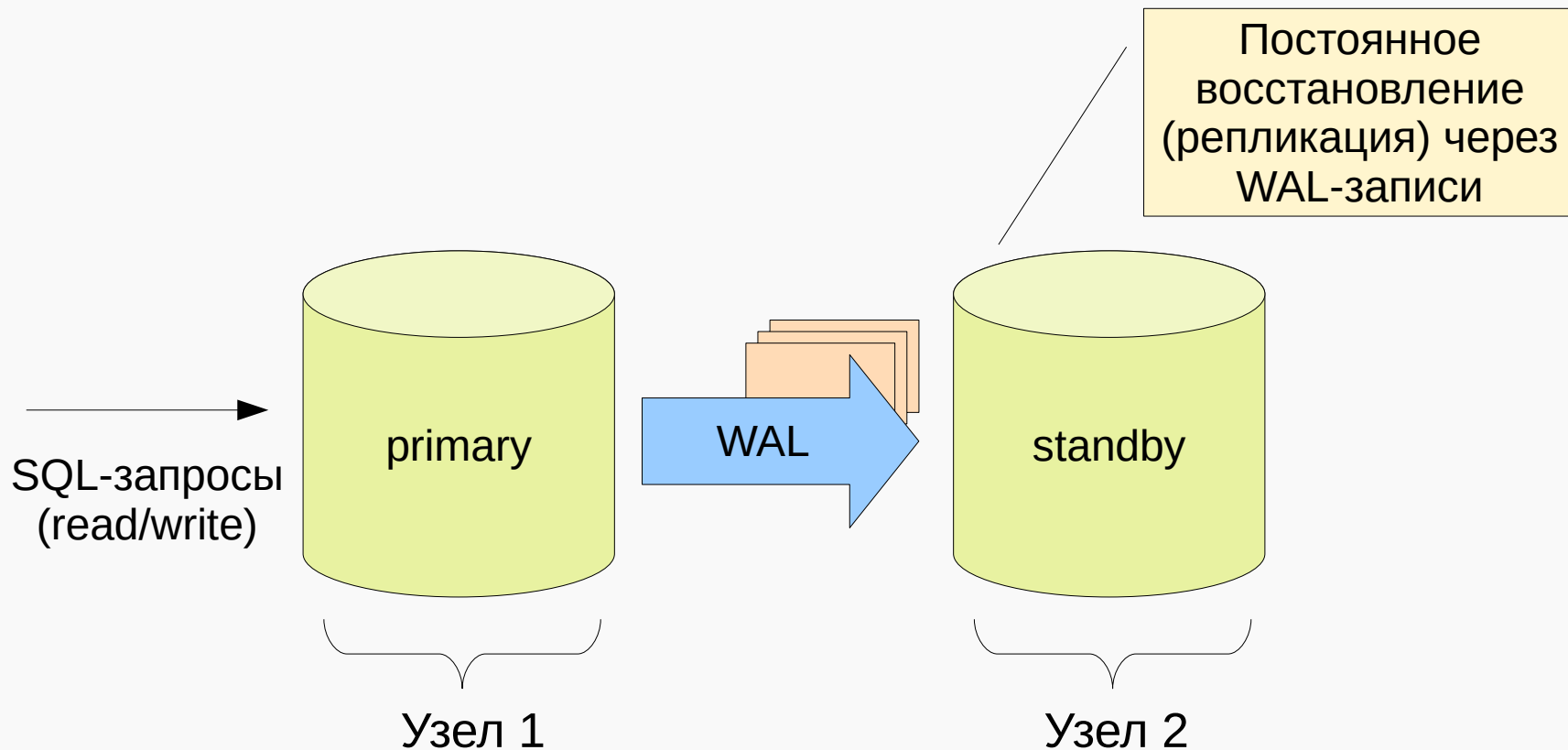
# Репликация в PostgreSQL

- В PostgreSQL репликация может быть:
  - **Физическая** — пересылаются файлы (WAL). standby — содержит копию primary:
    - ✓ **асинхронная** — данные для синхронизации состояний серверов БД пересылаются без подтверждения получения;
    - ✓ **синхронная** — пересылки данных идут с подтверждением.
  - **Логическая** — пересылаются декодированные из WAL операции.

# Физическая репликация

- **Физическая репликация** — пересылаются файлы (WAL). Standby — содержит копию primary.
- Требуется — `wal_level = replica` (по умолчанию).
- Поточковая репликация (**streaming replication**) — WAL-записи передаются по соединению между primary и standby:
  - Использует свой протокол для отправки WAL-файлов — WAL-sender, WAL-receiver.
- standby — осуществляет восстановление на основе полученных WAL для поддержки актуального состояния.

# Потоковая репликация



# Физическая репликация в PostgreSQL

- Пользователь для репликации:

```
CREATE USER PHYSREPL WITH REPLICATION ...
```

- pg\_hba.conf: host replication physrepl all md5
- backup: pg\_basebackup -d 'conn params' -D 'path'

} primary

- standby.signal в PGDATA.

- postgresql.conf:

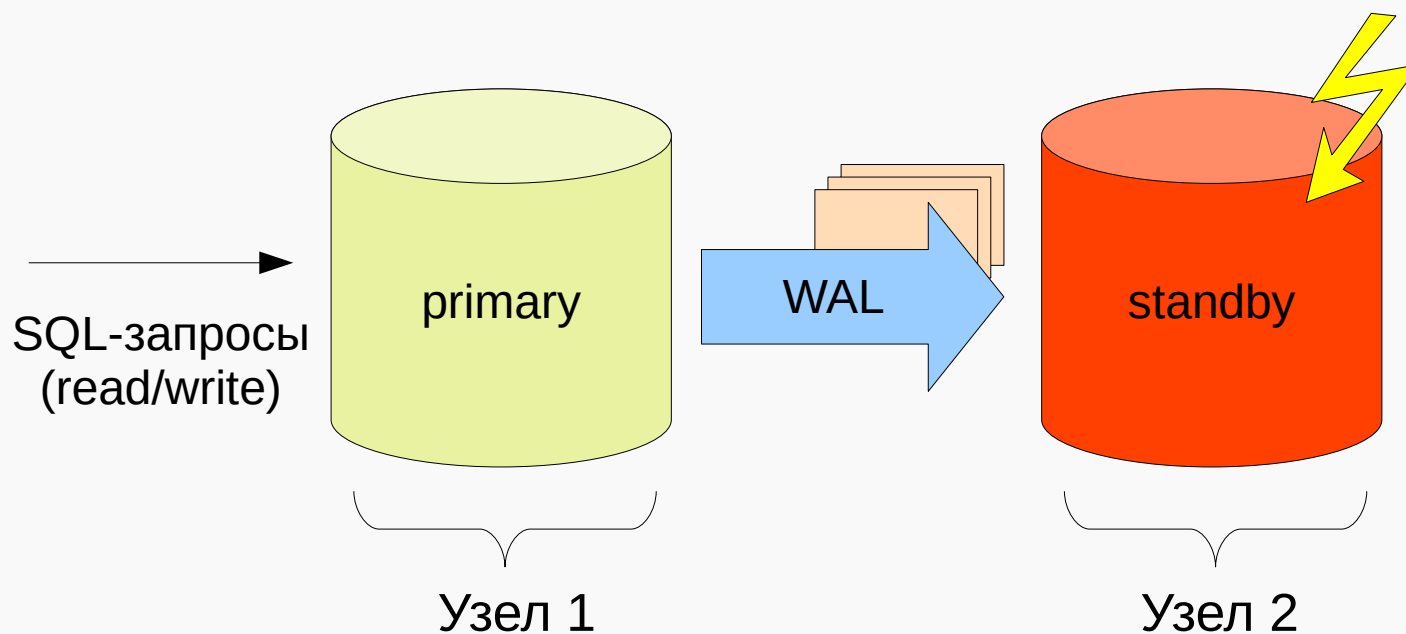
```
primary_conninfo = 'host=... user=physrepl'
```

} standby

# standby.signal

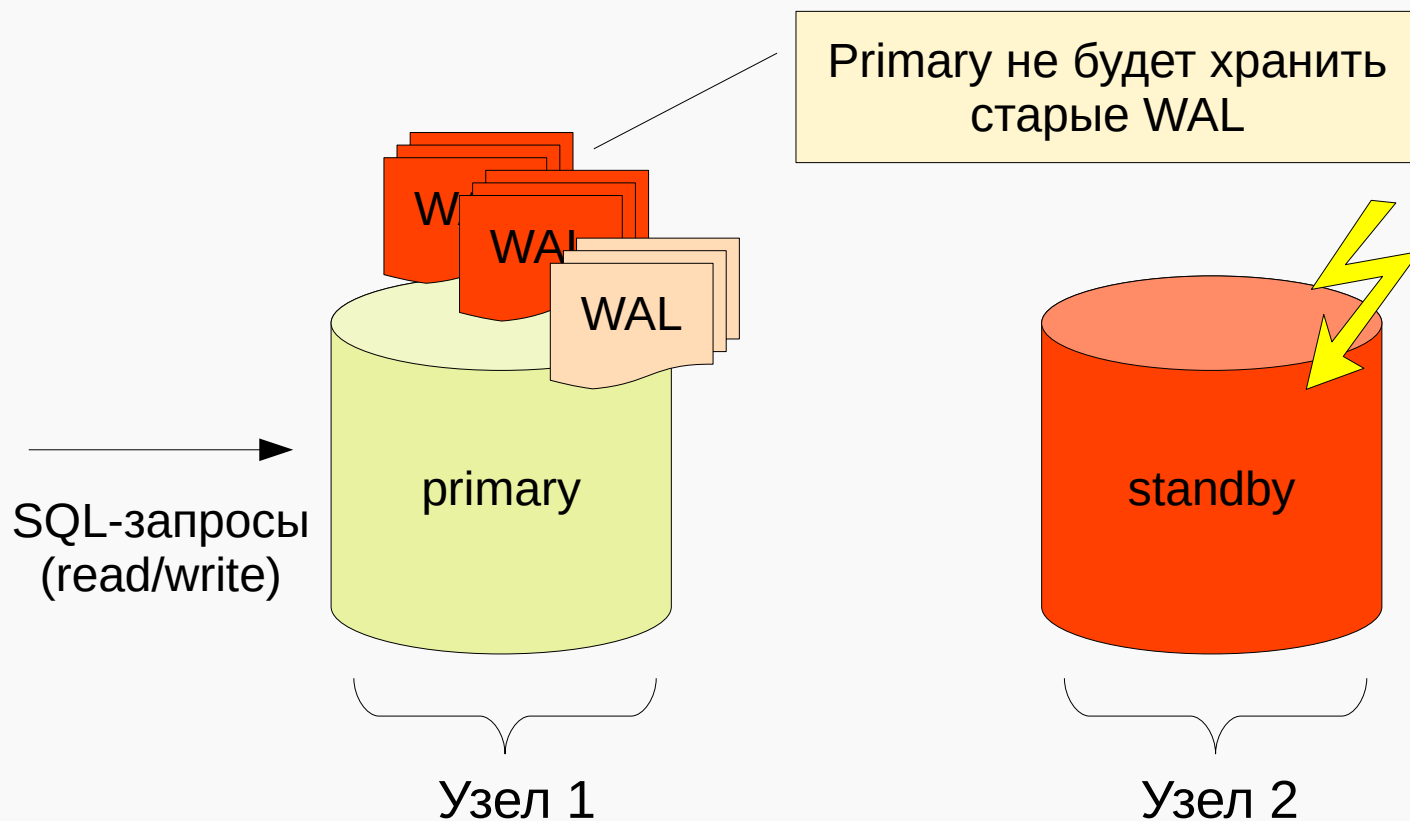
- Файл в PGDATA.
- Позволяет запустить сервер в режиме standby.
- Если установлен:
  - После обработки архивированных в данный момент WAL, сервер продолжает пытаться получить новые WAL-сегменты через:
    - ✓ `restore_command`
    - ✓ WAL-записи с мастера, используя `primary_conninfo`
- `primary_conninfo` — используется только, если `standby.signal` установлен.

# Что если произошла временная проблема со standby?

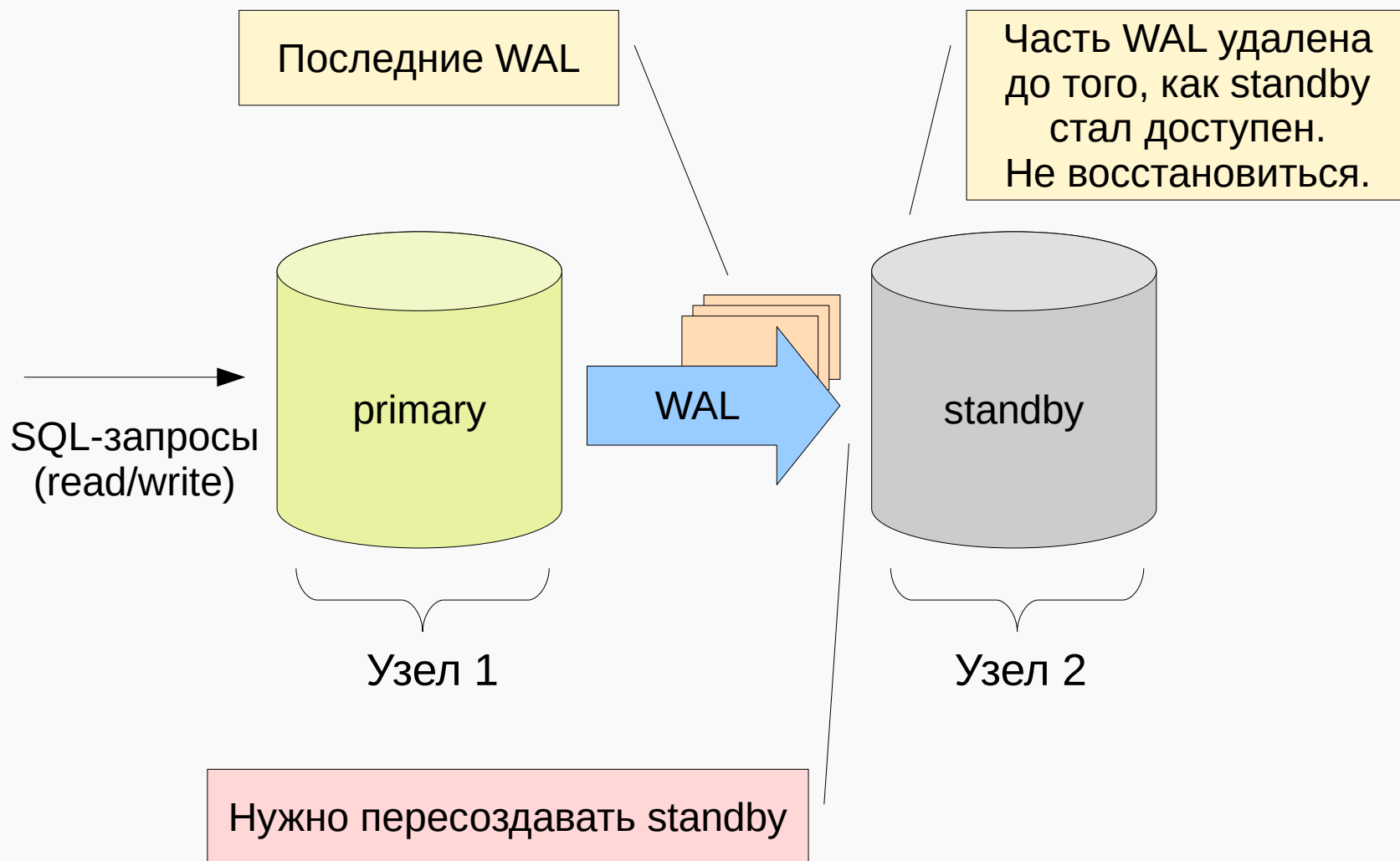




# Что если произошла временная проблема со standby?



# Что если произошла временная проблема со standby?



# wal\_keep\_size

- Задается в postgresql.conf.
- Определяет минимальное количество ненужных (для мастера) WAL-сегментов, которые должны храниться.
- По умолчанию — 0.
- Эти WAL-сегменты будут доступны standby:
  - **НЕТ** гарантии, что для standby будет достаточно этих сегментов.

- Позволяет хранить сегменты, пока они не отправлены standby.
- Таких сегментов может много накопиться (если standby не доступен).
- **max\_slot\_wal\_keep\_size** — позволяет ограничить число хранимых сегментов.
- Для использования:
  - нужно создать слот: `SELECT * FROM pg_create_physical_replication_slot('node2_slot');`
  - на standby: **primary\_slot\_name** = 'node2\_slot'.

# Полезные параметры/команды

- `--max_wal_senders` — число standby + 1. Если используется `--wal_method=stream` — еще +1 — для backup.
  - Определяет число процессов для отправки WAL (WAL sender).
- `--wal_method=stream` — при использовании `pg_basebackup` будет открыто второе соединение для передачи WAL параллельно с созданием копии.
  - Используется по умолчанию.
- Просмотр состояния реплики:
  - представление `pg_stat_replication`;

# Синхронная репликация

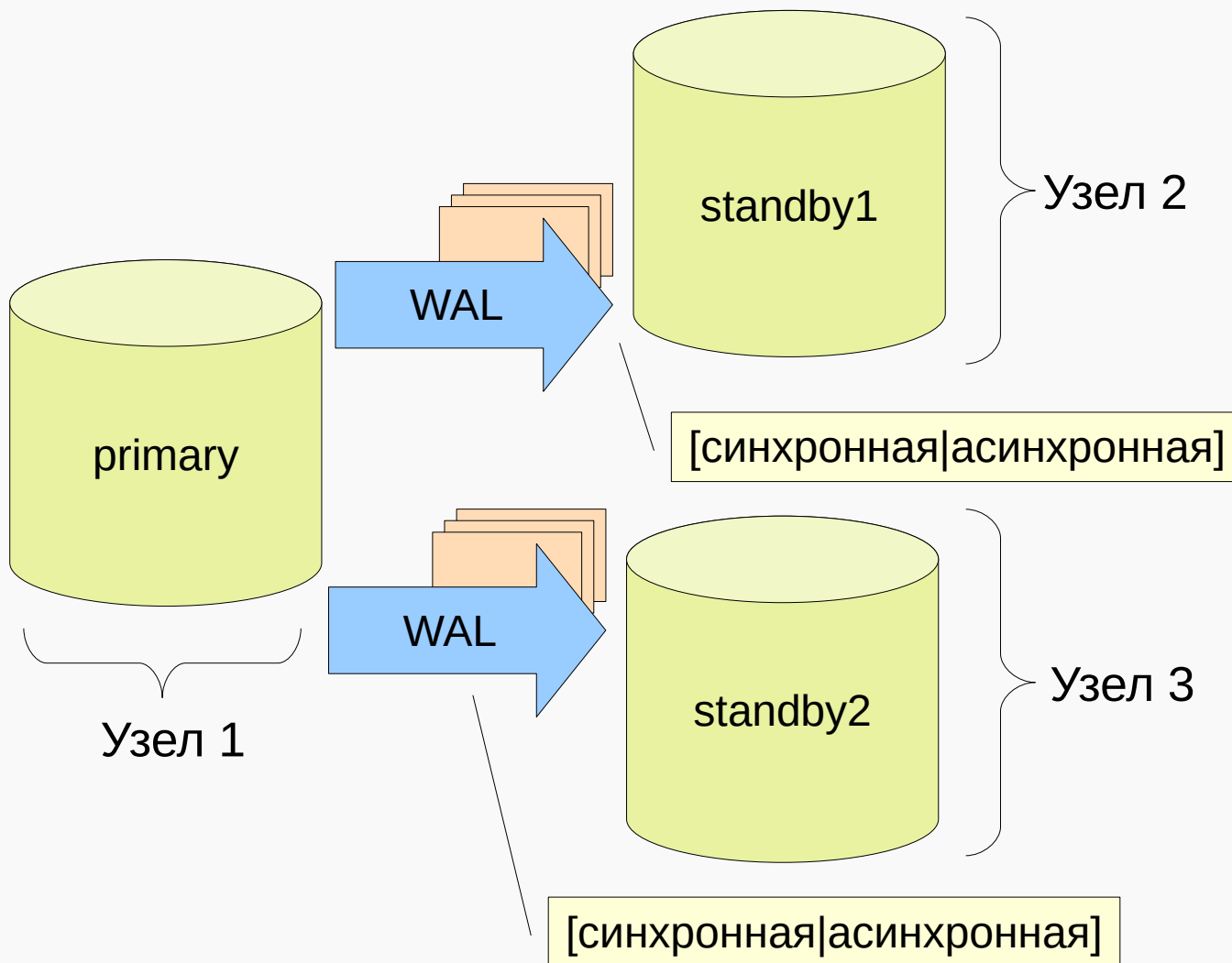
- Определяет, нужно ли ждать обработки WAL на репликах, когда происходит COMMIT транзакции.
- параметр `synchronous_commit` (default = on) + `synchronous_standby_names`:
  - по умолчанию включена асинхронная репликация.
- Дает возможность подтвердить, что транзакция реплицирована хотя бы на одном standby.
- Указать список standby серверов с синхронной репликацией:

```
synchronous_standby_names = 'node2, node3, node4'
```

# Synchronous commit

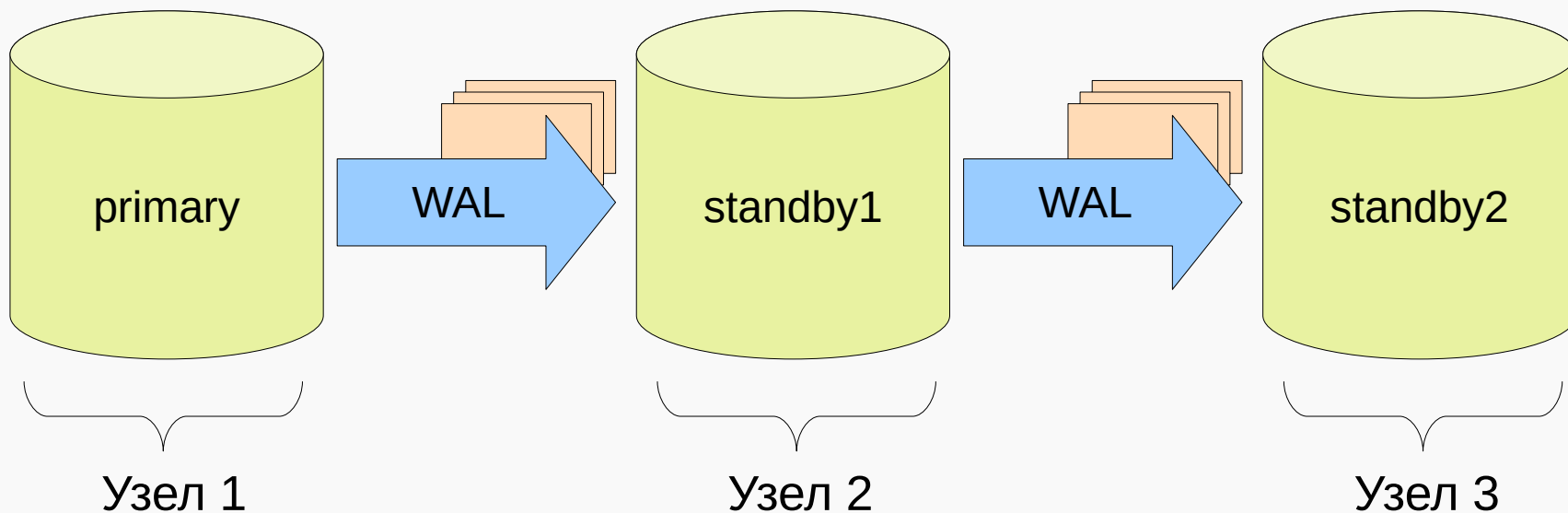
	Долгов-ть локального коммита	Долгов-ть сбой БД	Долгов-ть сбой ОС	standby - согл-ть запросов
remote_apply	+	+	+	+
on	+	+	+	
remote_write	+	+		
local	+			
off				

# Репликация на несколько узлов





# Ступенчатая репликация: cascading replication



# Особенности standby

- **warm standby** — режим standby, при котором логи (WAL), архивированные primary-сервером, собираются standby и происходит непрерывное восстановление состояния БД по полученным WAL. Warm standby не позволяет осуществлять запросы.
- **hot standby** — сервер работает в режиме чтения — кроме восстановления по WAL, позволяет осуществлять запросы к данным (read-only).

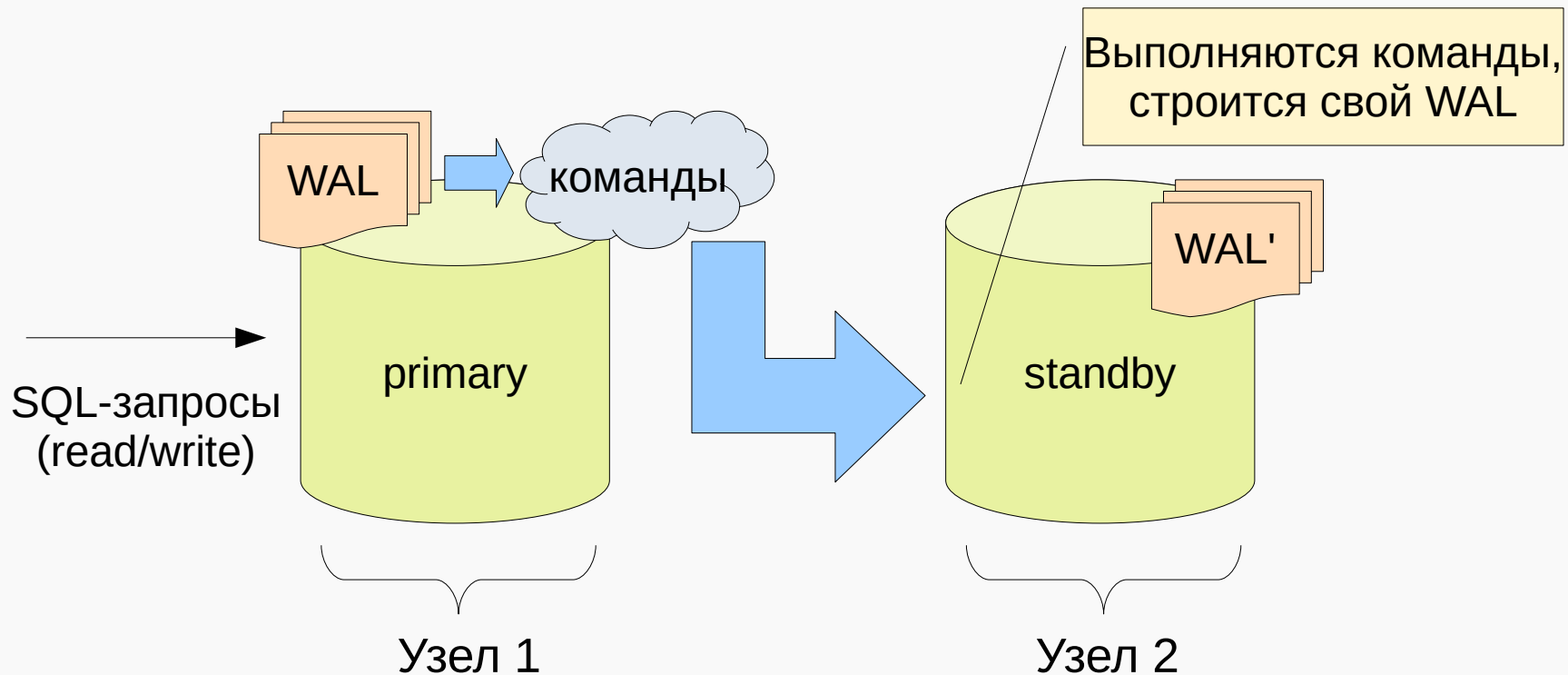
Параметр `hot_standby = true` (по умолчанию).

### 3. Логическая репликация данных в PostgreSQL

# Логическая репликация

- **Логическая репликация** — пересылаются команды, реконструированные из WAL.
- Требуется — `wal_level = logical`.
- `standby` — осуществляет восстановление на основе полученных команд для поддержки актуального состояния.
- Используется модель «публикаций» «подписок»:
  - после логического декодирования команды публикуются на мастере;
  - опубликованные команды могут получить подписанные `standby`.

# Логическая репликация



# Логическая репликация в PostgreSQL (master)

- wal\_level = logical
- Пользователь для репликации (master):  
`CREATE USER LOGICREPL WITH REPLICATION ...`
- pg\_hba.conf: host replication logicrepl all md5.
- Выдать права для LOGICREPL на нужные объекты.
- `CREATE PUBLICATION all_t_publ FOR ALL TABLES;`

primary

# Логическая репликация в PostgreSQL

- wal\_level = logical
- Воссоздать структуру таблиц.
- CREATE SUBSCRIPTION all\_t\_subscr  
CONNECTION 'user=logicropl ... dbname=db\_name'  
PUBLICATION all\_t\_publ;

standby

# Особенности

- Можно реплицировать часть БД (отдельный набор таблиц).
- Не поддерживается репликация DDL.
- Возможна репликация даже для узлов с разными версиями PostgreSQL.
- Standby может работать не в read-only режиме.



## Документация PostgreSQL:

<https://www.postgresql.org/docs/14/index.html>

## Лицензия PostgreSQL:

PostgreSQL is released under the PostgreSQL License, a liberal Open Source license, similar to the BSD or MIT licenses.

PostgreSQL Database Management System  
(formerly known as Postgres, then as Postgres95)

Portions Copyright © 1996-2022, The PostgreSQL Global Development Group

Portions Copyright © 1994, The Regents of the University of California

Permission to use, copy, modify, and distribute this software and its documentation for any purpose, without fee, and without a written agreement is hereby granted, provided that the above copyright notice and this paragraph and the following two paragraphs appear in all copies.

IN NO EVENT SHALL THE UNIVERSITY OF CALIFORNIA BE LIABLE TO ANY PARTY FOR DIRECT, INDIRECT, SPECIAL, INCIDENTAL, OR CONSEQUENTIAL DAMAGES, INCLUDING LOST PROFITS, ARISING OUT OF THE USE OF THIS SOFTWARE AND ITS DOCUMENTATION, EVEN IF THE UNIVERSITY OF CALIFORNIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

THE UNIVERSITY OF CALIFORNIA SPECIFICALLY DISCLAIMS ANY WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE SOFTWARE PROVIDED HEREUNDER IS ON AN "AS IS" BASIS, AND THE UNIVERSITY OF CALIFORNIA HAS NO OBLIGATIONS TO PROVIDE MAINTENANCE, SUPPORT, UPDATES, ENHANCEMENTS, OR MODIFICATIONS.

<https://www.interdb.jp/pg/index.html>