

PostgreSQL

The current, editable version of this book is available in Wikibooks, the open-content textbooks collection, at <https://en.wikibooks.org/wiki/PostgreSQL>

Permission is granted to copy, distribute, and/or modify this document under the terms of the [Creative Commons Attribution-ShareAlike 3.0 License](https://creativecommons.org/licenses/by-sa/3.0/).

Features

An outstanding Goal

The very first [concept paper](https://dsf.berkeley.edu/papers/ERL-M85-95.pdf) (https://dsf.berkeley.edu/papers/ERL-M85-95.pdf) - published in 1986 - defines a goal that distinguishes PostgreSQL from many other systems until today: "*provide user extendibility for data types, operators, and access methods*". This goal is reached. And it is not only available for users, even the internal implementation utilizes those interfaces to create system components.

Architecture

PostgreSQL implements a client/server model. Each *client process* connects to one *backend process* at the server site. Such backend processes are part of the *instance*, a group of many processes which act closely together and handle the data access. PostgreSQL does not use threading in the backend processes or elsewhere.

Features

Security

- Authentication methods: SCRAM-SHA-256, GSSAPI, SSPI, LDAP, RADIUS, Certificate, PAM
- Roles (users and groups) authorize access to data and execution of functions

Reliability

- Transactions with full ACID support and diverse isolation levels
- Savepoints (Sub-transactions)
- Multi-Version Concurrency Control (MVCC)
- Point-in-Time Recovery
- Partitioning of tables and indexes
- Synchronous, asynchronous, and logical replication
- Bi-Directional replication
- Publish/subscribe mechanism
- Parallel execution of single queries at multiple CPUs

Application Aspects

- Rich set of predefined data-types, i.a. JSON
- Support for arrays
- Composite type (similar to a *record* in some programming languages), constructor for rows via *row* keyword
- Check constraints
- Referential integrity with foreign keys
- Table inheritance
- Views, materialized views, updateable views

Extendability

- User defined data-types, operators, and index access methods
- User-defined functions, procedures, triggers, and procedural languages
- *Create extension* interface to create user-defined packages. Some publicly available examples:
 - *Foreign data wrappers* to other - PostgreSQL or non-PostgreSQL - databases or to the file-system
 - *PostGIS*: an extension for Spatial and Geographic Objects
 - *hstore*: a key/value storage

SQL Support

- High degree of conformance to the SQL standard: 170 out of 177 features
- Outer join, union, intersect, except

- Group by, grouping set, cube, rollup
- Common table expressions (CTE)
- Recursive queries, graph queries
- Window functions, analytic functions

Featured Platforms

PostgreSQL is available at diverse CPU architectures: x86, x86_64, IA64, PowerPC, PowerPC 64, S/390, S/390x, Sparc, Sparc 64, ARM, MIPS, MIPSEL, PA-RISC and runs on all major operating systems. ^[1]

- Linux
 - Red Hat family Linux (including CentOS/Fedora/Scientific/Oracle variants)
 - Debian GNU/Linux and derivatives
 - Ubuntu Linux and derivatives
 - SuSE and OpenSuSE
- macOS
- Windows (XP+)
- Solaris
- BSD
 - FreeBSD
 - OpenBSD
 - NetBSD
- AIX
- HP/UX

References

1. Operating Systems ^[1] (<https://www.postgresql.org/docs/current/static/supported-platforms.html>)

Download and Installation

Before you download PostgreSQL you must make two crucial decisions. First, decide whether to compile and install PostgreSQL from source code or to install from prebuilt binaries. Second (if you want to use any binary), you must know for which operating system you need the software. PostgreSQL supports most UNIX-based systems (including macOS) as well as Windows.

After you have made those decisions you can download and use the complete source code, an installer, a Bitnami Infrastructure Stack, or the pure binaries.

Start at the Source Code Level

The source code is available as a single packed file ^[1] or in a git repository ^[2]. To install from source you must download it to your local computer and compile it with a C compiler (at least C99-compliant, in most cases people use GCC (<https://gcc.gnu.org/>)) to the binary format of your computer. Details of the requirements ^[3], the download process, and the compilation steps ^[4] are available in the PostgreSQL documentation.

The advantages of working with the source code are that you can read and study it, modify it, or compile it on an exotic platform. But you must have some pre-knowledge and experience in handling specific tasks of your operating system, e.g.: working in a shell, installing additional programs, ...

The PostgreSQL documentation describes all details of the installation from source in the chapters:

- [Installation from Source Code on Unix \(https://www.postgresql.org/docs/current/installation.html\)](https://www.postgresql.org/docs/current/installation.html)
- [Installation from Source Code on Windows \(https://www.postgresql.org/docs/current/install-windows.html\)](https://www.postgresql.org/docs/current/install-windows.html)

Start with the Help of a Prebuild Program

In opposite to start at the source code level, it is relatively easy to use one of the pre-build programs or scripts. This is the preferred way for beginners. You can choose from several options:

- **Installer ^[5]:** This is the most comfortable way to download and install PostgreSQL on your local computer. The installer guides you not only through the installation steps, but also offers the option to install helpful additional tools and drivers. Installers are not available for all versions of all operating systems.
- **Bitnami infrastructure stack ^[6]:** Such stacks (WAPP, MAPP, LAPP, and others) offer the complete infrastructure (PostgreSQL, Apache Web Server, PHP) to run Web applications on Windows, macOS, or Linux.
- **Pure binaries ^[7]:** This is a listing of operating-specific commands which leads you thru the download and installation process of binaries.

Examples

Install binaries for Linux (Ubuntu) "PostgreSQL Apt Repository". <https://www.postgresql.org/download/linux/ubuntu/>.

```
# Create the file repository configuration:
sudo sh -c 'echo "deb http://apt.postgresql.org/pub/repos/apt $(lsb_release -cs)-pgdg main" >
/etc/apt/sources.list.d/pgdg.list'

# Import the repository signing key:
wget --quiet -O - https://www.postgresql.org/media/keys/ACCC4CF8.asc | sudo apt-key add -

# Update the package lists:
sudo apt-get update

# Install the latest version of PostgreSQL.
# If you want a specific version, use 'postgresql-12' or similar instead of 'postgresql':
sudo apt-get -y install PostgreSQL
```

Starting and stopping

```
sudo /etc/init.d/postgresql start
sudo /etc/init.d/postgresql stop
```

Windows

By default, PostgreSQL launches at each reboot so it may consume many resources. To avoid that, just execute *services.msc* and change the PostgreSQL service to start manually. Then, create a file *postgresql.cmd* containing:

```
net start postgresql-x64-9.5
pause
net stop postgresql-x64-9.5
```

As long as this script is launched as an administrator, the cluster with all its databases is available. Just press a key to shutdown the service.

More Information

The PostgreSQL [wiki](https://wiki.postgresql.org/wiki/Detailed_installation_guides) (https://wiki.postgresql.org/wiki/Detailed_installation_guides) offers a lot more information and hints about the installation steps.

After a successful installation, you will have

- The PostgreSQL binaries on your disc.
- A first cluster called *data* on your disc. The cluster consists of an empty database called *postgres* (plus two template databases) and a user resp. role called *postgres* as well.
- A set of Unix programs or a Windows service running on your computer. These programs/services handle the cluster with all its databases.

By default, PostgreSQL listens to port 5432. Possibly you must configure your firewall to reflect this situation.

Connect to the Database

After a successful installation, you have a cluster *data*, a database *postgres*, the database superuser *postgres*, and a new operating system user *postgres*. Login at the operating system level with the new operating system user. In a shell you can connect to the new database via the often used program `psql`. `psql` is a line-mode program similar to a shell and allows you to send SQL commands to the database.

```
$ # Example in Unix syntax
$ su - postgres
Password:
$
$ # psql --help      to see a detailed explanation of psql's options
$ # psql [OPTION]... [DBNAME [USERNAME]]
$ psql postgres postgres
psql (11.9 (Ubuntu 11.9-1.pgdg18.04+1))
Type "help" for help.

postgres=#
postgres=# \q -- terminate psql with backslash q or ctrl-d
$
```

The default prompt (prefix of every new line) of `psql` is `postgres=#`. After you have successfully started it, you can use SQL commands to communicate with the database. Here is an example that creates a new database user with the name 'nancy' - and deletes it afterward.

```
postgres=# CREATE USER nancy WITH ENCRYPTED PASSWORD 'ab8sxx5F4';
CREATE ROLE
postgres=#
postgres=# DROP USER nancy; -- delete the user
DROP ROLE
postgres=#
```

The database responds to every SQL command indicating its successful execution or an error. In the previous example `CREATE ROLE` means that the user is created.

Separation of Concerns

Please recap what you have so far: a cluster *data*, a database *postgres*, a user *postgres*. Furthermore, PostgreSQL divides every database into logical units which are called *schema*. Most objects reside in such a schema. The default schema is named *public* and exists in every database. The same applies to some special schemas where system information is stored. As far as you don't explicitly use schema-names, the schema *public* is utilized by default. This means that a `CREATE TABLE t (column_1 INTEGER);` command will create the table *t* in schema *public*.

We recommend avoiding the schema *public* for your data. Because *public* exists in every database, some tools use it to store their data there. Create and work in your own schema to have a clear distinction between system-, tools-, and user-data.

Second, avoid working with user *postgres*. This user account has very strong privileges and you should rarely use it. Create a user who acts as the stakeholder for your data, views, functions, trigger, etc. .

The following script creates a new user and its schema.

```
$ # start 'psql' as the original 'postgres' user with its strong privileges
$ psql postgres postgres
postgres=# -- the owner of the new schema shall be 'finance_master'
postgres=# CREATE USER finance_master WITH CREATEROLE LOGIN ENCRYPTED PASSWORD 'xxx';
CREATE ROLE
postgres=# -- the new schema 'finance' for your data
postgres=# CREATE SCHEMA finance AUTHORIZATION finance_master;
CREATE SCHEMA
postgres=# -- change 'search_path' (description of search_path: see below)
postgres=# ALTER ROLE finance_master SET search_path = finance, public;
ALTER ROLE
postgres=# \q
```

Start *psql* with the new user *finance_master*. We want him to work in schema *finance*, but every connection between *psql* and PostgreSQL acts at the database-level. It's not possible to specify an individual schema for a connection. Therefore PostgreSQL has implemented a mechanism called *search_path*. It simplifies the switching between schemas. *search_path* contains a list of schema names. Whenever you omit a schema name, this list is consulted to decide which schema to use. For our user *finance_manager* we have defined in the above ALTER ROLE command that he shall work in schema *finance* and - if there is no hit for his SQL command e.g. for a SELECT - the schema *public* is consulted next.

```
$ # -- first parameter of psql: database second parameter: user nothing for schema
$ psql postgres finance_master
postgres=# -- create a table
postgres=# CREATE TABLE t1 (column_1 INTEGER); -- table will be in schema 'finance' because of the 'search_path'
definition
CREATE TABLE
postgres=# -- you can use the schema name explicitly
postgres=# CREATE TABLE finance.t2 (column_1 INTEGER); -- table will be in schema 'finance' as well
CREATE TABLE
postgres=# -- it's possible to overwrite 'search_path' by using the schema name explicitly
postgres=# CREATE TABLE public.t3 (column_1 INTEGER); -- table will be in schema 'public'
CREATE TABLE
postgres=#
postgres=# \d -- this command lists schema, table, and owner names
List of relations
Schema | Name | Type | Owner
-----+-----+-----+-----
finance | t1 | table | finance_master
finance | t2 | table | finance_master
public | t3 | table | finance_master
postgres=#
```

References

1. Source code in a single packed file (via FTP) [2] (<https://www.postgresql.org/ftp/source/>)
2. Source code in a git repository [3] (<https://www.postgresql.org/docs/current/git.html>)
3. Requirements for compilation (Unix)[4] (<https://www.postgresql.org/docs/current/static/install-requirements.html>)
4. Installing from Source [5] (<https://www.postgresql.org/docs/current/installation.html>)
5. Installer [6] (<https://www.enterprisedb.com/downloads/postgres-postgresql-downloads>)
6. Bitnami stacks [7] (<https://bitnami.com/tag/postgresql>)
7. Download binaries [8] (<https://www.postgresql.org/download/>)

Managing the Instance

The PostgreSQL instance consists of several processes running on a server platform. They work together in a coordinated manner using common configuration files and a common start/stop procedure. Thus all are running or none of them.

The program `pg_ctl` controls and observes them as a whole. When you are logged in as user `postgres` you can start it from a shell. The simplified syntax is:

```
pg_ctl [ status | start | stop | restart | reload | init ] [-U username] [-P password] [--help]
```

status

When `pg_ctl` runs in the `status` mode, it lists the actual status of the instance.

```
$ pg_ctl status
pg_ctl: server is running (PID: 864)
/usr/lib/postgresql/9.4/bin/postgres "-D" "/var/lib/postgresql/9.4/main" "-c"
"config_file=/etc/postgresql/9.4/main/postgresql.conf"
$
```

You can observe, whether the instance is running or not, the process id (PID) of the postmaster, the directory of the cluster and the name of the configuration file.

start

When `pg_ctl` runs in the `start` mode, it tries to start the instance.

```
$ pg_ctl start
...
database system is ready to accept connections
$
```

When you see the above message everything works fine.

stop

When `pg_ctl` runs in the `stop` mode, it tries to stop the instance.

```
$ pg_ctl stop
...
database system is shut down
$
```

When you see the above message the instance is shut down, all connections to client applications are closed and no new applications can reach the database. The `stop` mode knows three different modes for shutting down the instance:

- *Smart* mode waits for all active clients to disconnect.
- *Fast* mode (the default) does not wait for clients to disconnect. All active transactions are rolled back and clients are forcibly disconnected.
- *Immediate* mode aborts all server processes immediately, without a clean shutdown.

Syntax: `pg_ctl stop [-m s[mart] | f[ast] | i[mmediate]]`

restart

When `pg_ctl` runs in the `restart` mode, it performs the same actions as in a sequence of `stop` and `start`.

reload

In the `reload` mode the instance reads and reloads its configuration file.

init

In the `init` mode the instance creates a complete new cluster with the 3 databases *template0*, *template1*, and *postgres*. This command needs the additional parameter `-D datadir` to know at which place in the file system it shall create the new cluster.

Automated start at boot time

In most cases it is necessary that PostgreSQL starts immediately after the server boots. Whether this happens - or not - may be configured in the file `start.conf`, which is located in the special directory `$PGDATA` (Debian/Ubuntu) or in the main directory of the cluster (RedHat). There is only one entry and its allowed values are:

- `auto`: automatically start/stop
- `manual`: do not start/stop automatically, but allow manually managing as described above
- `disabled`: do not allow manual startup with `pg_ctlcluster` (this can be easily circumvented and is only meant to be a small protection for accidents)

Tools

Tools for database administration (DBA) tasks such as backups, restores, and cleanups are mostly not part of the SQL standard. Vendor specific database products usually include a combination of database specific tools and SQL extensions for administration purposes. PostgreSQL provides a set of both PostgreSQL specific tools and SQL extensions. The main ones are described here as well as some reliable external tools.

psql

psql is a client program which is delivered as an integral part of the PostgreSQL downloads. Similar to a bash shell it is a line-mode program and may run on the server hardware or a client. *psql* knows two kinds of commands:

- Commands starting with a backslash, eg: `\dt` to list tables. Those commands are interpreted by *psql* itself.
- All other commands are sent to the instance and interpreted there, e.g.: `SELECT * FROM mytable;`

Thus it is an ideal tool for interactive and batch SQL processing. The whole range of PostgreSQL SQL syntax can be used to perform everything that can be expressed in SQL.

```
-----
$ # start psql from a bash shell for database 'postgres' and user 'postgres'
$ psql postgres postgres
postgres=#
postgres=# -- a standard SQL command
postgres=# CREATE TABLE t1 (id integer, col_1 text);
CREATE TABLE
postgres=# -- display information about the new table
postgres=# \dt t1
      List of relations
 Schema | Name | Type  | Owner
-----+-----+-----+-----
 public | t1   | table | postgres
(1 row)
postgres=#
postgres=# -- perform a PostgreSQL specific task - as an example of a typically DBA action
postgres=# SELECT pg_start_backup('pitr');
pg_start_backup
-----
0/2000028
(1 row)
postgres=#
postgres=# -- terminate psql
postgres=# \q
$
-----
```

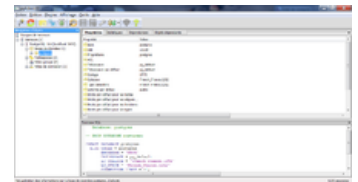
Here are some more examples of *psql* 'backslash'-commands

- `\h` lists syntax of SQL commands
- `\h SQL-command` lists syntax of the named SQL-command
- `\?` help to all 'backslash' commands
- `\l` lists all databases in the current cluster
- `\echo :DBNAME` lists the current database (consider upper case letters). In most cases the name is part of the *psql* prompt.
- `\dn` lists all schemas in the current database
- `\d` lists all tables, views, sequences, materialized views, and foreign tables in the current schema
- `\dt` lists all tables in the current schema
- `\d+ TABLENAME` lists all columns and indexes in table TABLENAME
- `\du` lists all users in the current cluster
- `\dp` lists access rights (privileges) to tables, ...
- `\dx` lists installed extensions
- `\o FILENAME` redirects following output to FILENAME
- `\t` changes output to 'pure' data (no header, ...)
- `\! COMMAND` executes COMMAND in a shell (outside *psql*)
- `\q` terminates *psql*

pgAdmin

pgAdmin is a tool with a graphical user interface for Unix, Mac OSX and Windows operating systems. In most cases it runs on a different hardware than the instance. For the major operating systems it is an integral part of the download, but it is possible to download the tool separately (<http://www.pgadmin.org/download/>).

pgAdmin significantly extends the functionalities of *psql* with intuitive, graphical representations of database objects, eg. schemas, tables, columns, users, result lists, query execution plans, dependencies between database objects, and much more. To give you a first impression of the surface, some screenshots (<http://www.pgadmin.org/screenshots/>) are online.



pgAdmin

Since 2016 *pgAdmin 3* is superseded by *pgAdmin 4*. *pgAdmin 4* is a complete re-implementation - written as a web application in Python. You can run it either on a web server using a browser, or standalone on a workstation.

phpPgAdmin

phpPgAdmin is a graphical tool which offers features that are similar to those of *pgAdmin*. It is written in PHP, therefore you additionally need Apache and PHP packages.

phpPgAdmin is not part of the standard PostgreSQL downloads. It is distributed via Sourceforge (<http://sourceforge.net/projects/phpPgAdmin/>). The project was largely dormant over the past few years, but as of 2019 has been updated with support for PHP7 and PostgreSQL 12^[1].



PhpPgAdmin

Other Tools

There are a lot of other general tools (https://wiki.postgresql.org/wiki/Community_Guide_to_PostgreSQL_GUI_Tools) and monitoring tools (https://wiki.postgresql.org/wiki/Monitoring#Postgres-centric_monitoring_solutions) with a GUI interface. Their functionality varies greatly from pure SQL support up to entity-relationship and UML support. Some of the tools are open/free source, others proprietary.

References

1. https://github.com/phpPgAdmin/phpPgAdmin/releases/tag/REL_7-12-0

Configuration

The main configuration file is *postgresql.conf*. It is divided into several sections according to different tasks. The second important configuration file is *pg_hba.conf*, where authentication definitions are stored.

Both files reside in the special directory \$PGDATA (Debian/Ubuntu) or in the main directory of the cluster (RedHat).

Some configuration items have a dynamic nature, and will take effect with a simple `pg_ctl reload`. Others require a restart of the instance `pg_ctl restart`. The comments in the default configuration files state which of the two actions has to be taken.

postgresql.conf

File Locations

The value of *data_directory* defines the location of the cluster's main directory. In the same way the value of *hba_file* defines the location and the name of the above mentioned *pg_hba.conf* file (host based authentication file), where rules for authentication are stored - some more details are shown [below](#).

Connections

In the connections section you define the port number (default: 5432), with which client applications can reach the instance. In addition the maximum number of connections is defined as well as SSL, IP and TCP settings.

Resources

The main definition in the resources section is the size of shared buffers. It determines, how much space is reserved to "mirror" the content of data files within PostgreSQL's buffers in RAM. The predefined default value of 128 MB is relatively low.

Secondly, there are definitions for the work and the maintenance memory. They determine the RAM sizes for sorts, create index commands, These two RAM areas exist per connection and are used individually by them whereas the shared buffers exist only once for the whole instance and are used concurrently by multiple processes.

Additionally there are some definitions concerning *vacuum* and *background writer* processes.

WAL

In the WAL section there are definitions for the behaviour of the WAL mechanism.

First, you define a WAL level out of the four possibilities *minimal*, *archive*, *hot_standby*, and *logical*. Depending on the decision, which kind of archiving or replication you want to use, the WAL mechanism will either write only basic information to the WAL files or include additional information. *minimal* is the basic method which is always required for every crash recovery. *archive* is necessary for any archiving action, which includes the point-in-time-recovery (PITR) mechanism. *hot_standby* adds information required to run read-only queries on a standby server. *logical* adds information necessary to support logical decoding.

Additionally and in correlation to the WAL level *archive* there are definitions which describe the archive behaviour. Especially the 'archive_command' is essential. It contains a command which copies WAL files to an archive location.

Replication

If you use replication to a different server, you can define the necessary values for master and standby server in this section. The master reads and pays attention only on the master-definitions and the standby only to the standby-definitions (you can copy this section of 'postgresql.conf' directly from master to standby). You must define the WAL level to an appropriate value.

Tuning

The tuning section defines the relative costs of different operations: sequential disc I/O, random disc I/O, process one row, process one index entry, process one function-call or arithmetic operation, size of effective RAM pages (PostgreSQL + OS) per process which will be available at runtime. These values are used by the query planner during its search for an optimal query execution plan. The values are not real values in sense of milliseconds or number of CPU cycles. They are only a rough guideline for the query planer and relative to each other. The real values are calculated during the query execution may differ significantly.

There is also a subsection concerning costs for the genetic query optimizer, which - in opposite to the standard query optimizer - implements a heuristic searching for optimal plans.

Error Logging

The error logging section defines the amount, location and format of log messages which are reported in error situations or for debugging purposes.

Statistics

In the statistics section you can define - among other things - the amount of statistic collection for parsing, planing and execution of queries.

pg_hba.conf

The *pg_hba.conf* file (host-based authentication) contains rules for client access to the instance. All connection attempts of clients which do not satisfy these rules are rejected. The rules restrict the connection type, client IP address, database within the cluster, user-name, and authentication method.

There are two main connection types: local connections (*local*) via sockets and connections via TCP/IP (*host*). The term *local* refers to the situation, where a client program resides on the same machine as the instance. The client may override the *local* connection and use the *host* connection type by using the TCP/IP address syntax (e.g.: 'localhost:5432') of the cluster.

The client IP address is a single IPv4 or IPv6 address or a masking of a net-segment via a CIDR mask.

The database and the client user name must be given explicitly or may be abbreviated by the key word "ALL".

There are different authentication methods

- *trust*: don't ask for any password
- *reject*: don't allow any access
- *password*: ask for a password
- *md5*: same as 'password', but the transfer of the password occurs MD5-encrypted
- *peer*: trust the client, if he uses the same database username as his operation system username (only applicable for local connections)

Since the *pg_hba.conf* records are examined sequentially for each connection attempt, the order of the records is significant. The first match between defined criteria and properties of incoming connection requests hits.

BackupAndRecovery

Overview

Creating backups is an essential task for every database administrator. If the hardware crashes or there is software corruption, the DBA must ensure that a database can be restored with minimal data loss. PostgreSQL offers multiple strategies to support the DBA in achieving this goal.

Backup technology can be divided into two classes: cold backups and hot backups. A cold backup is a backup taken when no database file is open. In the case of PostgreSQL this means that the instance is stopped, the backup is taken and the instance is restarted. A hot backup is a backup taken during normal working operations. Applications can perform read and write actions in parallel with the backup creation. There are different types of hot backups. The main types are: logical, physical, and physical with PITR. These will be described in more detail in the section on hot backups.



Cold Backup (Offline Backup)

A cold backup is a backup taken when the PostgreSQL instance is not running and includes a **consistent** copy of all files which constitute the cluster with all of its databases. Cold backups are also called *offline* backups.

There is only one way to create a consistent useful cold backup: the PostgreSQL instance should be stopped by issuing the `pg_ctl stop` command. This will disconnect all applications from all the cluster's databases. After the instance is shut down, make a backup using one of the standard operating system utilities (cp, tar, dd, rsync, etc.) to create a copy of all cluster files to a secure location, such as: on the disk of a different server, on a backup system at a SAN or intranet, a tape system or other reliable location.

To be a successful backup, the backup must include the following cluster files:

- Include all files under the directory node where the cluster is installed. The environment variable `$PGDATA` points to this directory and usually resolves to something like, `.../postgresql/9.4/main`. Use `echo $PGDATA` on the command-line or in `psql`, `show data_directory;` to retrieve the directory node path.
- Include configuration files, which may be in `$PGDATA`, but can also be located elsewhere, for example, as they are on the default Ubuntu install. These are the main configuration files: `postgresql.conf`, `pg_hba.conf`, and `pg_ident.conf`. Their locations can be found by running the following commands from the `psql` utility:

```
show config_file;
show hba_file;
show ident_file;
```

- Include all tablespace files. These files are located elsewhere on the file-system. Their locations can be found by looking at the symlinks in the `$PGDATA/pg_tblspc` directory:

```
cd $PGDATA/pg_tblspc
ls -lt
```

Caution! One may try to backup only special parts of a cluster, eg. a huge file which represents a table on a separate partition or tablespace - or the opposite: everything except the huge file. Even if the instance is shut down during the generation of such a partial copy, copies of this kind are useless. The restore of a cold backup needs all data files and meta-information files of the cluster to re-create the cluster.

Advantages

- A cold backup is easy to do and to restore.

Disadvantages

- A continuous 7x24 operation mode of any of the databases in the cluster is not possible.

- It is not possible to backup smaller parts of a cluster like a single database or table.
- Partial restores are not possible. Restores must include all the cluster files.
- After a crash, any data changes that occur after the most recent cold backup are lost. Only the data in the backup will be restored.

How to Recover

It is strongly recommended to do the following steps on a test host to verify the restore before doing them on the production server.

In the case of a crash, restore the data from a cold backup by performing the following steps:

- Stop the instance.
- Backup the original files of the crashed cluster: the files in the \$PGDATA path and the configuration files.
- Delete all original files of the crashed cluster: the files in the \$PGDATA path and the configuration files.
- Copy the files of the backup to their original places.
- Start the instance. It should start in the normal way, without any special message.

Hot Backup (Online Backup)

In contrast to *cold backups*, *hot backups* are taken while the instance is running and applications may change data as the backup is taken.

Hot backups are also known as *online backups*.

Logical Backup

A logical backup is a **consistent** copy of the data within a database or some of its parts. These backups are created with the utilities `pg_dump` and `pg_dumpall`. The instance must be running for these tools to operate. While they may be run in parallel to applications, they create consistent snapshots as of the time they are called. For example, if an application changes some data values during the creation of a backup, the backup takes the old values whereas the application sees the new values.

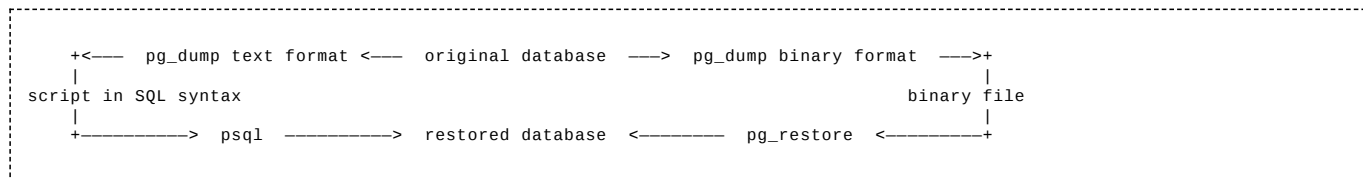
Logical backups run in serializable transactions. This is possible because of PostgreSQL's MVCC (Multi-version concurrency control) implementation.

pg_dump

`pg_dump` works at the database level and can backup specific subsets of a database such as individual tables. It is able to dump data, schema definitions or both. The parameters `--data-only` and `--schema-only` restrict the output to information respective to the given flag.

`pg_dump` supports two output formats; [text format](#) and [binary format](#). The format type is chosen by `pg_dumps` parameter `--format`. The [text format](#) contains SQL commands like `CREATE` and `INSERT`. Files created in this format may be used by `psql` to restore the backed-up data. The [binary format](#) is also called the *archive format*. To restore files with this format you must use the `pg_restore` tool.

The following diagram visualizes the cooperation of `pg_dump`, `psql` and `pg_restore`.



Some Examples:

```

$ # dump complete database 'finance' in text format to a file
$ pg_dump --dbname=finance --username=boss --file=finance.sql
$
$ # restore database content (to a different or an empty database)
$ psql --dbname=finance_x --username=boss <finance.sql
$
$
$
$ # dump table 'person' of database 'finance' in binary format to a file
$ pg_dump --dbname=finance --username=boss --table=person --format=c --file=finance_person.archive
$
$ restore table 'person' from binary archive
$ pg_restore --dbname=finance_x --username=boss --format=c <finance_person.archive
$

```

pg_dumpall

The `pg_dumpall` command works at the cluster level and can backup up important cluster level objects like user/roles and their rights. `pg_dumpall` without detailed parameters will dump the complete contents of the cluster: all data, and structures of all databases plus all user definitions and definitions of their rights. The parameter `--globals-only` can be used to restrict its behavior to dump global objects only. `pg_dumpall` outputs in the [text format](#).

Advantages

- Continuous 7x24 operation mode is possible.
- Small parts of cluster or database may be backup-ed or restored.
- When you use the text format you can switch from one PostgreSQL version to another or from one hardware platform to another.

Disadvantages

- The text format uses much space, but it compresses well.

How to Recover

As mentioned in the above diagram the recovery process depends on the format of the dump. Text files are in standard SQL syntax. To recreate objects from this format you can use SQL utilities like `psql`. Binary files are in the *archive format*. They can only be used by the utility `pg_restore`.

Physical Backup

A physical backup is an **inconsistent** copy of the files of a cluster, created with operating system utilities like `cp` or `tar` taken at a time whereas applications modify data. At first glance such a backup seems to be useless. To understand its purpose, you must know PostgreSQL's recover-from-crash strategy.

At all times and independent from any backup/recovery action, PostgreSQL maintains *Write Ahead Log (WAL)* files - primarily for crash-safety purposes. Such *WAL files* contain *log records*, which reflects all changes made to the data and the schema. Prior to their transfer to the data files of the database the *log records* are stored in the (sequentially written) *WAL file*. In the case of a system crash those *log records* are used to recover the cluster to a consistent state during restart. The recover process searches the timestamp of the last *checkpoint* and replays all subsequent *log records* in chronological order against this version of the cluster. Through that action the cluster gets recovered to a consistent state and contains all changes up to the last COMMIT.

The existence of a physical backup, which is inconsistent by definition but contains its *WAL files*, in combination with this recovery-from-crash technique can be used for backup/recovery purposes. To implement this, you have to restore the previous taken physical backup (including its *WAL files*). When the instance starts again, it uses the described recovery-from-crash technique and replays all *log records* in the *WAL files* against the database files. In the exact same way as before, the cluster comes to a consistent state and contains all changes up to the point in time when the backup-taken has started.

Please keep in mind, that physical backups work only on cluster level, not on any finer granularity like database or table.

Physical backups without PITR sometimes are called *standalone hot physical backup*.

Advantages

- Continuous 7x24 operation mode is possible.

Disadvantages

- Physical backup works only on cluster level, not on any finer granularity like database or table.
- Without PITR (see below) you will lose all data changes between the time, when the physical backup is taken, and the crash.

How to Take the Backup and Recover

To use this technique it is necessary to configure some parameters in the `postgres.conf` file for *WAL* and archive actions. As the usual technique is *Physical Backup plus PITR* we describe it in the next chapter.

Physical Backup plus PITR

The term *PITR* stands for *Point In Time Recovery* and denotes a technique, where you can restore the cluster to any point in time between the creation of the backup and the crash.

The *Physical Backup plus PITR* strategy takes a physical backup plus all WAL files, which are created since the time of taking this backup. To implement it, three actions must be taken:

- Define all necessary parameters in *postgres.conf*
- Generate a physical backup
- Archive all arising WAL files

If a recovery becomes necessary, you have to delete all files in the cluster, recreate the cluster by copying the physical backup to its original location, create the file *recovery.conf* with some recovery-information (especially to what location WAL files have been archived) and restart the instance. The instance will recreate the cluster according to its parameters in *postgres.conf* and *recovery.conf* to a consistent state including all data changes up to the last COMMIT.

Advantages

- Continuous 7x24 operation mode is possible.
- Recover with minimal data loss.
- Generating WAL files is the basis for additional features like *replication*.

Disadvantages

- Physical backup works only on cluster level, not on any finer granularity like database or table.
- If your database is very busy and changes a lot of data, many WAL files may arise.

How to Take the Backup

Step 1

You have to define some parameters in *postgres.conf* so that: WAL files are on the level 'archive' or higher, archiving of WAL files is activated and a copy command is defined to transfers WAL files to a fail-safe location.

```
# collect enough information in WAL files
wal_level = 'archive'
# activate ARCHIVE mode
archive_mode = on
# supply a command to transfer WAL files to a failsafe location (cp, scp, rsync, ...)
# %p is the pathname including filename. %f is the filename only.
archive_command = 'scp %p dba@archive_server:/postgres/wal_archive/%f'
```

After the parameters are defined, you must restart the cluster `pg_ctl restart`. The cluster will continuously generate WAL files in its subdirectory *pg_wal* (*pg_xlog* in Postgres version 9.x and older) in concordance with data changes in the database. When it has filled a WAL file and must switch to the next one, it will copy the old one to the defined archive location.

Step 2

You must create a *physical* or *base backup* with an operating system utility during the instance is in a special 'backup' mode. In this mode the instance will perform a checkpoint and create some additional files.

```
$ # start psql and set the instance to 'backup' mode, where it creates a checkpoint
$ psql -c "SELECT pg_start_backup('AnyBackupLabel');"
$
$ # copy the cluster's files
$ scp -r $PGDATA dba@archive_server:/postgres/whole_cluster/
$
$ # start psql again and finish 'backup' mode
$ psql -c "SELECT pg_stop_backup();"
$
```

If you like to do so, you can replace the three steps by a single call to the utility *pg_basebackup*.

Step 3

That's all. All other activities are taken by the instance, especially the continuous copy of completely filled WAL files to the archive location.

How to Recover

To perform a recovery the original *physical* or *base backup* is copied back and the instance is configured to perform recovery during its start.

- Stop the instance - if it is still running.
- Create a copy of the crashed cluster - if you have enough disc space. Maybe, you will need it in a later stage.
- Delete all files of the crashed cluster.
- Recreate the cluster files from the base backup.
- Create a file *recovery.conf* in `$PGDATA`. It must contain a command similar to: `restore_command = 'scp dba@archive_server:/postgres/wal_archive/%f %p'`. This copy command is the reverse of the command in *postgres.conf*, which saved the WAL files to the archive location.
- Start the instance. During startup the instance will copy and process all WAL files found in the archive location.

The fact, that *recovery.conf* exists, signals the instance to perform a recovery. After a successful recovery this file is renamed.

If you want to recover to some previous point in time prior to the occurrence of the crash (but behind the creation of the backup), you can do so by specifying this point in time in the *recovery.conf* file. In this case the recovery process will stop before processing all archived WAL files. This feature is the origin of the term *Point-In-Time-Recovery*.

In summary the *recovery.conf* file may look like this:

```
restore_command      = 'scp dba@archive_server:/postgres/wal_archive/%f %p'  
recovery_target_time = '2016-01-31 06:00:00 CET'
```

Additional Tools

There is an open source project *Barman* ^[1], which simplifies the steps of backup and recovery. The system helps you, if you have to manage a lot of servers and instances and it becomes complicate to configure and remember all the details about your server landscape.

References

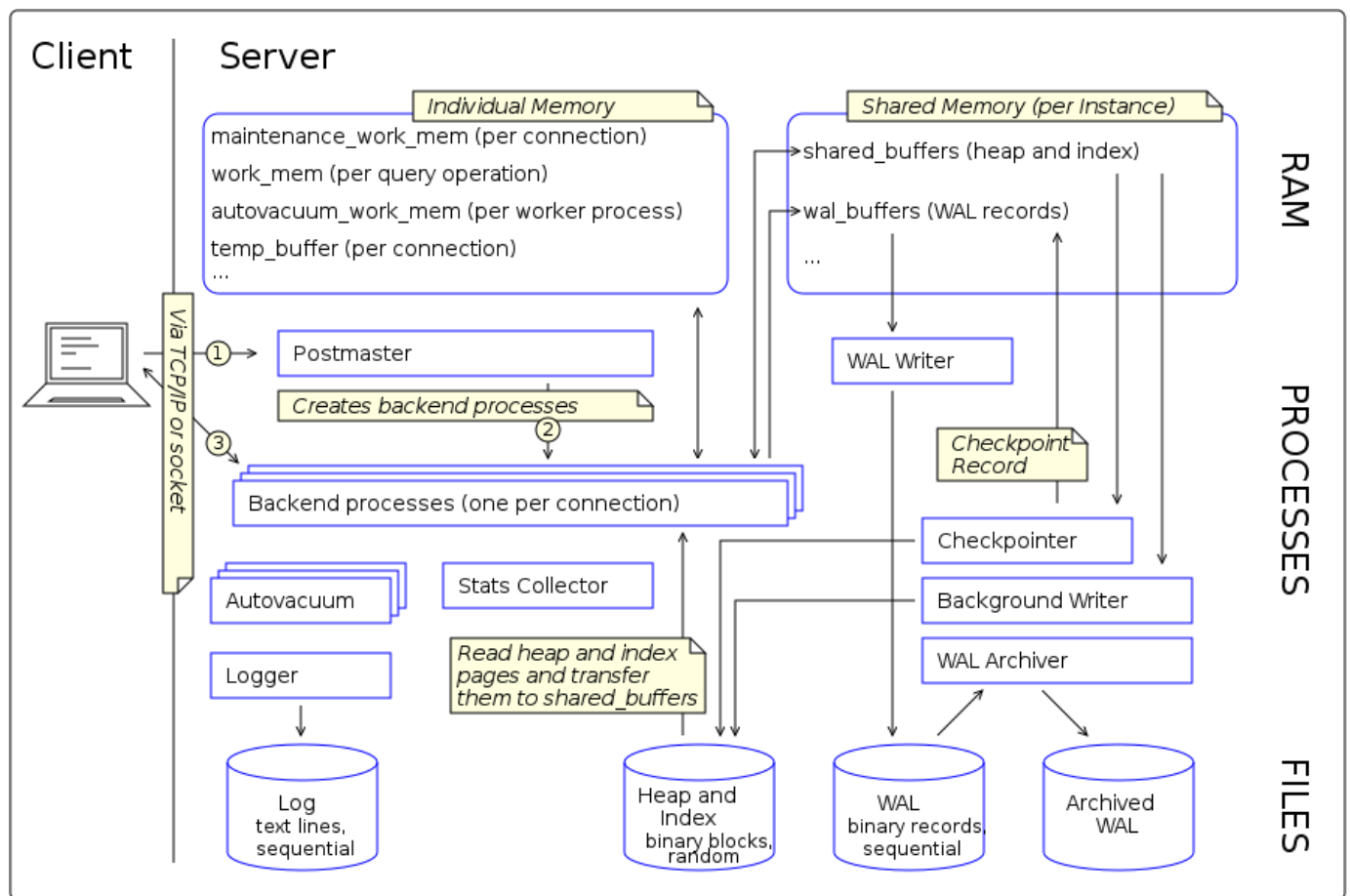
1. Barman [9] (<http://www.pgbarman.org/>)

Architecture

PostgreSQL implements a client-server architecture. Each `Client` process connects to one `Backend` process at the server site.

Clients do not have direct access to database files and the data stored in them. Instead, they send requests to the `Server` and receive the requested data from there. The server launches a single process for each client connection. Such a `Backend` process handles the client's requests by acting on the `Shared Memory`. This leads to other activities (file access, WAL, vacuum, ...) of the `Instance`. The `Instance` is a group of server-side processes acting on the common `Shared Memory`. PostgreSQL does not use threading.

At startup time, the `Instance` is launched by a single process denoted the `Postmaster`. It loads configuration files, allocates `Shared Memory`, and starts the other collaborating processes of the `Instance`: `Background Writer`, `Checkpointer`, `WAL Writer`, `WAL Archiver`, `Autovacuum`, `Statistics Collector`, `Logger`, and more. Subsequently, the `Postmaster` listens to its configured system port. In response to new client connection attempts he launches new `Backend` processes and delegates authentication, communication, and the handling of all further requests to them. The next figure visualizes the main aspects of `RAM`, processes, files, and their collaboration.



Client requests like `SELECT` or `UPDATE` usually lead to the necessity to read or write data. This is carried out by the client's `Backend` process. Such `I/O`-activities are **not done directly on the disks**. Instead, they are done in a cache in the `Shared Memory` that mirrors the file pages. Accesses to such caches are much faster than accesses to disk. Read accesses affect only the cache whereas write accesses are accomplished by writing to a log, the so-called write-ahead-log or `WAL`.

`Shared Memory` is limited in size and it can become necessary to evict pages. As long as the content of such pages hasn't changed, this is not a problem. But they may be modified. Modified pages are called `dirty pages` (or `dirty buffers`) and before they can be evicted they must be written back to disk. The `Background Writer` processes and the `Checkpointer` takes care of that. They ensure that the cache is - after a short time delay - in sync with files. The synchronization from `RAM` to disk consists of two steps.

First, whenever the content of a page changes, a `WAL record` is created containing the delta-information (difference between the old and new content) and stored in another area of `Shared Memory`. During a `COMMIT` or earlier the `WAL Writer` process reads them and appends them to the end of the current `WAL file`. Such sequential writes are faster than writes to random positions of heap and index files. All `WAL records` created from one `dirty page` must be transferred to disk before the `dirty page` itself can be transferred to disk in the second step.

Second, the Background Writer process transfers dirty buffers from Shared Memory to files. Because I/O activities can block other processes, it starts periodically and acts only for a short period. Doing so, its extensive - and expensive - I/O activities are spread over time, avoiding debilitating I/O peaks. The Checkpointer process also transfers dirty buffers to file.

The Checkpointer process creates a **Checkpoint** by writing and flushing all older dirty buffers, all older WAL records, and finally a special Checkpoint record to disk. Therefore a Checkpoint is a point in the sequence of transactions at which it is guaranteed that the heap and index files have been updated with all dirty pages before that Checkpoint.

WAL files contain the changes made to the data. Such 'delta information' is used in the case of a system crash for recovery (database backup + WAL files --> database immediately before the crash). Hence, WAL files shall be duplicated and preserved at a safe place until the next database backup is taken. This is the duty of the WAL Archiver process. He can be configured to run a script that copies WAL files, as soon as they are full and a switch to the next one takes place. Of course, such copies should be done to a separate disk or server for security reasons. But it's also a good idea to store the original WAL files on a different disk than heap and index files. Such a separation boosts performance. It can be done using a symbolic link pointing from the original WAL directory to a directory at a different disk.

The Autovacuum process marks old versions of records in the heap and index files that are no longer used by any transaction as 'finally deleted'. Hence it releases the space occupied by them for reuse. The need for such a process results from the MVCC architecture.

The Statistics Collector collects counters about accesses to SQL objects like tables, rows, indexes, pages, ... and stores them in system tables.

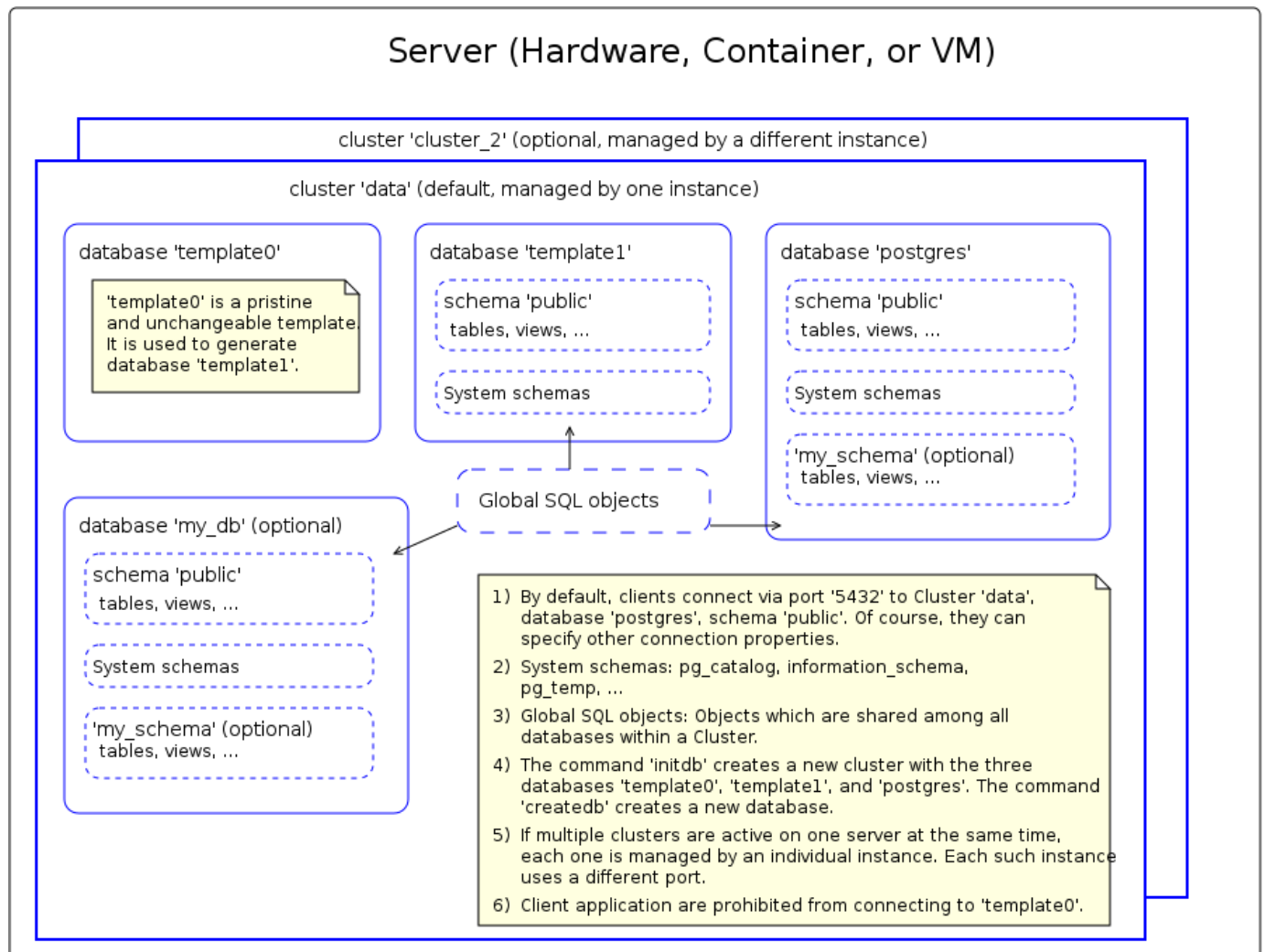
The Logger writes text lines about more or less serious events that may happen during database accesses, e.g., wrong password, no permission, long-running queries, etc. to a sequential file.

Architecture Cluster

Overview

A *Server*, which is some hardware, a container, or a VM, contains one or more Database Clusters (*Cluster* for short). Every cluster is controlled by exactly one instance. If there are many clusters and instances on the same server, the ports of the instances must differ from each other as well as the root directories of the clusters.

Each newly created cluster contains the three databases *template0*, *template1*, and *postgres*, each of the three databases contain the schema *public* as well as the system schemas *pg_catalog*, *information_schema*, *pg_temp*, and some more. Tables, views, and most other SQL objects reside in such schemas. DBAs can create more clusters, databases, schemas, or SQL objects.



Initialization Phase

Clusters are created with the command `initdb`. *template0* is the very first database during the creation phase of any cluster. In a second step, database *template1* is generated as a copy of *template0*, and finally database *postgres* is generated as a copy of *template1*. Later, the DBA can create more databases within that cluster, e.g.: *my_db*, with the command `createdb`. Just like at the beginning, the new database will be a copy of *template1*. Due to the unique role of *template0* as the pristine original of all other databases, no client is allowed to connect to it and modify it. But the DBA can change *template1*.

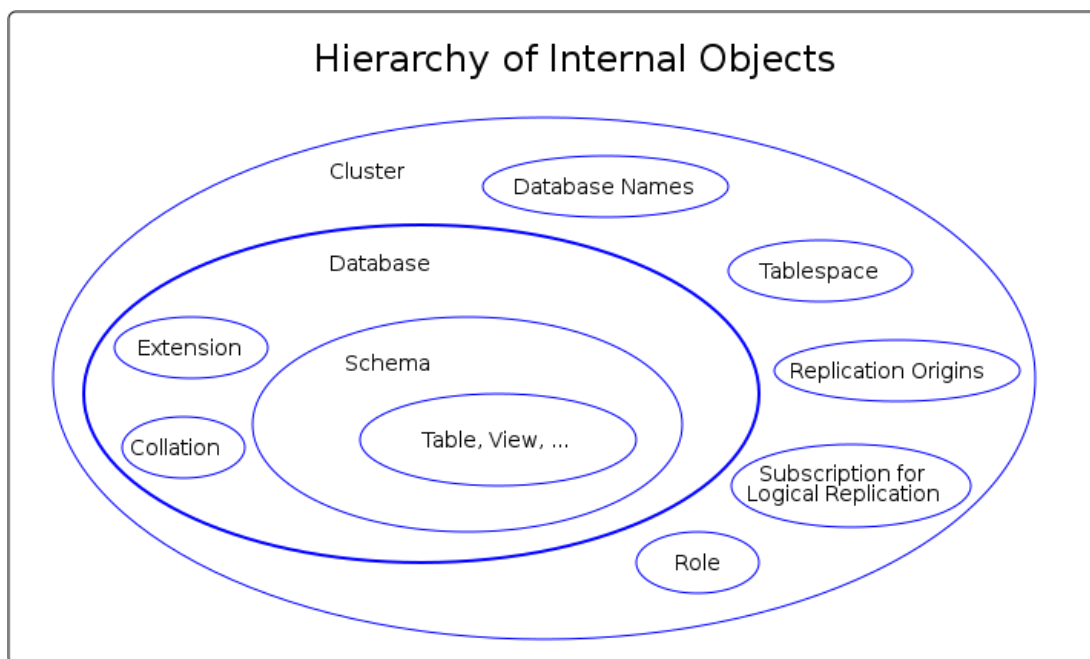
Connections

Client connections act at the database level and can access data and SQL objects within any schema of the connected database, as far as they are permitted to do so. If they need access to any object of a different database within the same or another cluster, special techniques like `foreign-data wrapper` (FDW) or `dblink` are required (or they use multiple connections and synchronize them at the client-side).

SQL objects

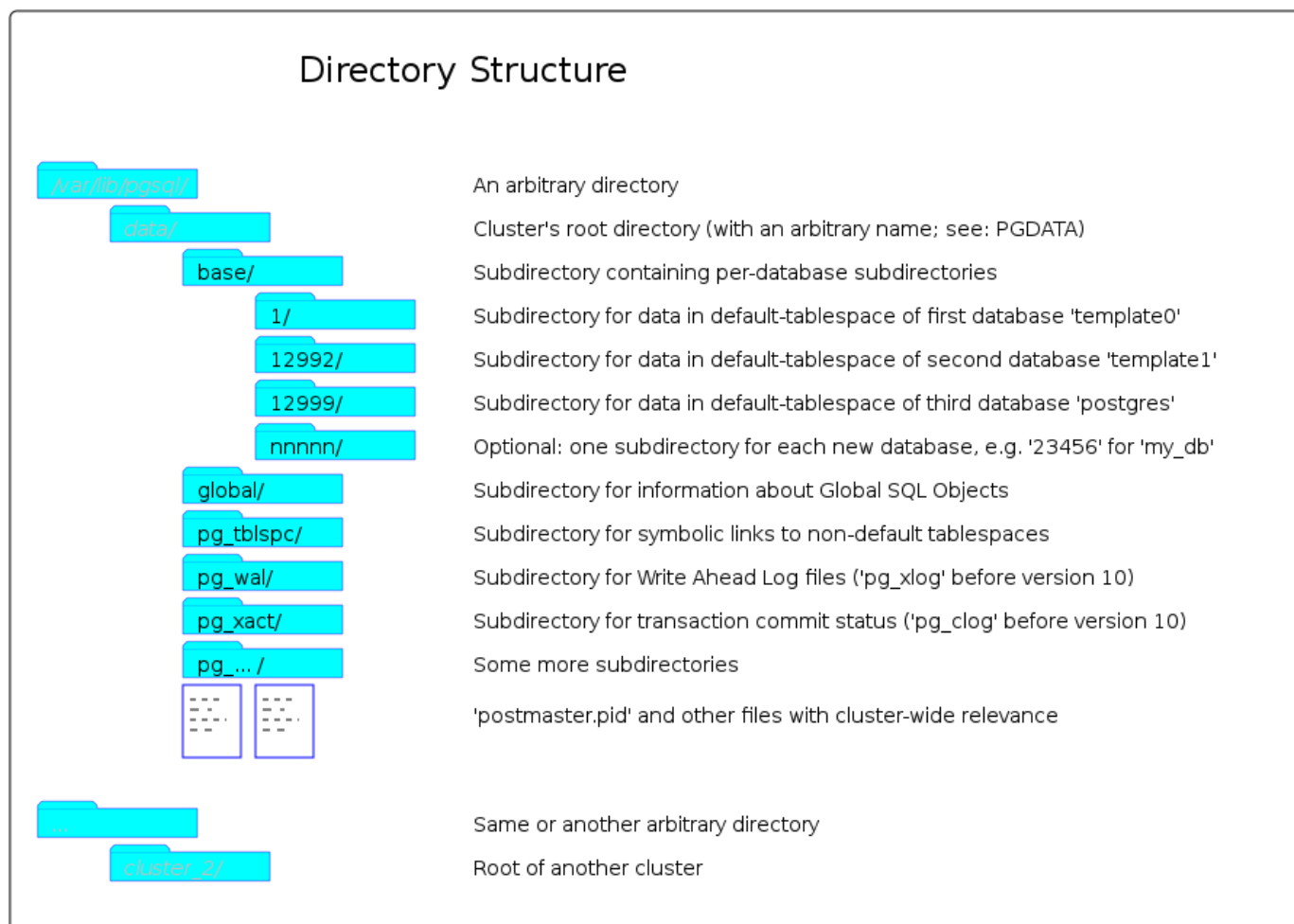
We use the term *SQL object* for all objects which you can create with the SQL command `CREATE . . .`, e.g.: database, schema, table, view, materialized view, index, constraint, sequence, function, procedure, trigger, role, data type, domain, operator, tablespace, extension, foreign-data wrapper, and much more. Such SQL objects are arranged in a hierarchical manner:

- Database names, tablespaces, and roles (users) are known at the cluster level. E.g.: As mentioned above, a connection works at the database level. Nevertheless, when you create a new role with such a connection, the role is also known by all other databases of the same cluster.
- Extensions, e.g.: PostGIS, reside at the database level. After installing an extension, all schemas of this database can use it. But within the other databases of the same cluster, the extension is not known.
- Schemas are part of a database. Some of them are predefined.
 - *pg_catalog* is a schema with tables that describe most of the SQL objects of that database, especially all tables and views. They even describe themselves. *information_schema* is a similar schema. It contains several tables and views of *pg_catalog* in a way that conforms to the SQL standard.
 - *public* acts as the default schema. It should not contain user-defined SQL objects. Instead, it is recommended to create one or more additional schemas to manage application-specific objects like tables or triggers. To access objects in such additional schemas they can be fully qualified, e.g. `my_schema.my_table`, or by changing the `search_path`.
- There are different types of SQL objects within a schema: 'relation'-like objects (table, view, materialized view, index, sequence, foreign-table), function, procedure, trigger, constraint, data type, domain, operator, and more.
 - SQL objects in one schema are different from SQL objects in different schemas, even if they use the same name, e.g.: table *t1* in *my_schema1* is different from *t1* in *my_schema2*.
 - The names of 'relation'-like objects, data types, and domains are unique within their schema: e.g.: you cannot have a table *employee* and a view *employee* in the same schema.



Architecture Directories

PostgreSQL organizes durable (persistent) data as well as volatile state information about transactions or replication actions in the file system. Every cluster has its root directory somewhere in the file system. In many cases, the environment variable `PGDATA` points to this directory. The following graphic uses `data`, which is the default, as the name of the cluster's root directory.



The cluster's root directory contains many subdirectories and some files, all of which are necessary to store durable as well as temporary information. The root's name can be selected as desired, but the names of its subdirectories and files are constant respectively determined by PostgreSQL. The following paragraphs describe the most important subdirectories and files.

`base` contains one subdirectory per database. The names of those subdirectories consist of numbers. These are the internal Object Identifiers (OID), which are numbers to identify their definition in the System Catalog.

Within the database-specific subdirectories of `base`, there are many files: one or more for every heap and index. Again, the filenames consist of numbers. Those files are accompanied by files for the Free Space Maps (suffixed `_fsm`) and Visibility Maps (suffixed `_vm`), which contain optimization information. An example for filenames is: `3083`, `3083_fsm`, `3083_vm`.

Another subdirectory is `global`. It contains files with information about SQL Objects which are not restricted to a certain schema, but known and relevant at the schema level.

In `pg_tblspc`, there are symbolic links that point to directories that are outside of the root directory tree, e.g. at a different disk. Heap and index files of non-default tablespaces reside there. Those defined within the default tablespace reside in the database-specific subdirectories.

The subdirectory `pg_wal` contains the WAL files. They arise and grow in parallel with data changes in the cluster and remain as long as they are required for recovery, archiving, or replication.

The subdirectory `pg_xact` contains information about the status of each transaction: `in_progress`, `committed`, `aborted`, or `sub_committed`.

In the root directory, there are some files. In many cases, the configuration files of the cluster are stored here. Also, if the instance is up and running, the file *postmaster.pid* exists here (by default, but other locations are possible). It contains the process ID (pid) of the Postmaster process which has started the instance and controls it.

Transactions

All data-changing operations like INSERT, UPDATE, or DELETE must run within a surrounding construct which is called a TRANSACTION. Transactions are created with the SQL command BEGIN and finished with either COMMIT or ROLLBACK. During the lifetime of the transaction the changes to the database are written only preliminarily. At the end, COMMIT finishes the transaction regularly and commits all intended data changes, or ROLLBACK aborts the transaction and reverts all those preliminary changes.

In addition to this explicit usage of SQL keywords to manage transactions, some of PostgreSQL's client libraries create implicitly a new transaction if one of the data-changing operations doesn't run in an explicitly created transaction. In this case, the operation is automatically committed immediately after its execution.

```
BEGIN; -- establish a new transaction
UPDATE accounts SET balance = balance - 100.00 WHERE name = 'Alice';
UPDATE accounts SET balance = balance + 100.00 WHERE name = 'Bob';
COMMIT; -- finish the transaction

-- this UPDATE runs as the only command of an implicitly created transaction ...
UPDATE accounts SET balance = balance - 100.00 WHERE name = 'Alice';

-- ... and this one runs in another transaction
UPDATE accounts SET balance = balance + 100.00 WHERE name = 'Bob';
```

Hint: When working within a procedure or function, there is a DECLARE ... BEGIN ... END; construct to define 'blocks'. In this context the meaning of BEGIN (no semikolon after BEGIN!) differs from what is explained here. You can avoid the ambiguity by using the keywords START TRANSACTION; as an alternative for BEGIN; in the context of transactions. Besides that, START TRANSACTION; conforms to the SQL standard.

Transactions generate a great relief for applications. Especially for business logic that must execute many statements as a consistent unit - like the above money transfer from one bank account to another -, there is no need to take individual actions after an error occurred in the middle of a transaction. In many cases it's enough to restart the transaction or to handle errors in a unified way.

Transactions in PostgreSQL guarantee that all requirements of the ACID paradigm are fulfilled, [see next chapter](#).

Sub-Transactions

Within a transaction the keyword SAVEPOINT defines and denotes a position, to which the transaction may be rolled-back.

```
-- The transaction will insert the values 1 and 3, but not 2.
BEGIN;
INSERT INTO my_table VALUES (1);
SAVEPOINT my_savepoint;
INSERT INTO my_table VALUES (2);
ROLLBACK TO SAVEPOINT my_savepoint;
INSERT INTO my_table VALUES (3);
COMMIT;
```

ACID

The ACID paradigm is a cornerstone of database management systems. With respect to data modifications, the paradigm demands that transactions must fulfill certain requirements and have to guarantee that they are satisfied not only during regular operations but also in all cases of minor and major problems like mutual-locking, connection-loss, server-down, disk-full, disk-crash,

Please note especially that the requirements are defined at the **transaction level**. They are named:

- Atomicity
- Consistency
- Isolation
- Durability

As shown in the previous chapter every single write operation or a collection of write operations is embedded in a transaction. Read operations may also be part of a transaction.

Atomicity

All writing operations within a transaction create a single, undividable unity. Either all of them succeed or none. Writing operations to different tables are an example of such a situation. Another example is the decrease of one person's bank account and the associated increase of another bank account during a money transfer.

Consistency

At the end of a transaction, the database is in a consistent state. All defined integrity rules like uniqueness, check constraints, foreign-key and primary-key definitions are fulfilled. Furthermore, all involved triggers have been successfully executed. It's possible that during the lifetime of a transaction those rules may be broken, e.g. the foreign key relationship of two nodes in a doubly-linked list.

In essence, a transaction transfers the database from one consistent state to another consistent state.

Isolation

In many cases, transactions run in parallel. But the database system gives them the illusion that they act one after the other. Depending on the chosen isolation level, the exact behavior of competing read and write operations may differ. Nevertheless, PostgreSQL guarantees in all cases, that read operations never block write operations and write operations never block read operations.

Durability

Durability guarantees that after a successful termination (COMMIT) of a transaction, the carried-out changes keep in the database, even if a significant problem like a disk crash occurs. PostgreSQL implements this by saving the data changes not only in the data files but - redundant - also in Write-Ahead-Log (WAL) files. Therefore it is recommended that data and WAL files shall be stored on different disks.

Visibility

Some exemplary Problems

It's obvious that every transaction 'sees' all data changes, it has been carried out during its lifetime, without problems. But there are situations where more than one process wants to read or write the same data during an overlapping time interval of their transactions or even at the same point in time, which is possible on servers with multiple CPUs or a disk array. In such cases, different types of conflicts and suspicious effects may occur.

Applications may or may not accept the effects resulting from such competing situations. They can choose different levels of *isolation* against the activities of other transactions depending on their needs. The level defines which effects they are willing to accept and which not. Higher levels mean that fewer effects can occur but the database system must work harder and that the overall throughput decreases.

Here are some examples with two transactions T_A and T_B . Both don't perform a COMMIT if not explicitly noted.

- T_A reads the row with $id = 1$. T_B reads the same row. T_A increases column X by 1. T_B increases the same column by 1. What will be the result? There is the danger of a 'Lost update'.
- T_A changes a value of the row with $id = 1$. What shall T_B see if it reads the same row? T_A may perform a ROLLBACK. (Uncommitted read)
- T_A reads the row with $id = 1$. T_B reads the same row, changes a value and performs a COMMIT. T_A reads the row again. In comparison to its first read, it will see a different value. (Non-repeatable read)
- T_A reads all rows with $status = 'ok'$. T_B inserts an additional row with $status = 'ok'$ and performs a COMMIT. T_A reads all rows with $status = 'ok'$ again and receives a different number of rows. (Phantom read)
- T_A reads and changes the row with $id = 1$. T_B reads and changes the row with $id = 2$. T_B wants to read and change the row with $id = 1$. Because T_A has not yet committed its changes, T_B must wait for T_A . T_A wants to read and change the row with $id = 2$. Because T_B has not yet committed its changes, T_A must wait for T_B . (Deadlock)

PostgreSQL's Solutions

The SQL standard describes the 3 effects (or problematic situations) 'Uncommitted read', 'Non-repeatable read', and 'Phantom read' and defines 4 levels of isolation between transactions: READ UNCOMMITTED, READ COMMITTED, REPEATABLE READ, and SERIALIZABLE. Every level is stricter than its predecessor and prevents more effects, which means e.g. that a 'Non-repeatable read' is possible in level READ COMMITTED but not in REPEATABLE READ or SERIALIZABLE.

PostgreSQL implements those levels (<https://www.postgresql.org/docs/current/transaction-iso.html>). But, as a consequence of its MVCC model, it implements some aspects a little stricter than they are demanded by the standard. If a transaction requests the level READ UNCOMMITTED, PostgreSQL handles it always as a READ COMMITTED, which leads to the overall behavior that all uncommitted changes are invisible to all other transactions at any level - only committed changes can be seen by other transactions.

Examples

The following examples act on a table $t1$ with the two columns id and col and a single row.

```
DROP TABLE IF EXISTS t1;
CREATE TABLE t1 (id INTEGER, col INTEGER);
INSERT INTO t1 VALUES (1, 100);
SELECT * FROM t1;
id | col
---+---
 1 | 100
(1 row)
```

Uncommitted read

The example shows that PostgreSQL solely shows committed rows to other transactions.

Transaction A

Transaction B

```
BEGIN;
SET TRANSACTION ISOLATION LEVEL READ UNCOMMITTED;
-- you can shorten the two commands into one:
-- BEGIN TRANSACTION ISOLATION LEVEL READ UNCOMMITTED;
```

```
UPDATE t1 SET col=101 WHERE id=1;
SELECT col FROM t1 WHERE id=1;
-- 101
```

```
-- 'READ UNCOMMITTED' acts equal to 'READ COMMITTED'
-- other transactions solely sees committed rows!
BEGIN TRANSACTION ISOLATION LEVEL READ UNCOMMITTED;
SELECT col FROM t1 WHERE id=1;
-- 100 (the committed one!)
```

```
COMMIT;
SELECT col FROM t1 WHERE id=1;
-- 101
```

```
SELECT col FROM t1 WHERE id=1;
-- 101 (again: the committed one!)
COMMIT; -- no real effect
SELECT col FROM t1 WHERE id=1;
-- 101
```

Lost update

The example shows that PostgreSQL prevents 'lost update' in the lowest level of isolation - as well as in all other levels. (The table *t1* contains its original values.)

Transaction A

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
SELECT col FROM t1 WHERE id=1;
-- 100
UPDATE t1 SET col=col+1 WHERE id=1;
SELECT col FROM t1 WHERE id=1;
-- 101
```

Transaction B

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
SELECT col FROM t1 WHERE id=1;
-- 100
UPDATE t1 SET col=col+1 WHERE id=1;
-- UPDATE is queued and must wait for the
-- COMMIT of transaction A
.
```

```
COMMIT;
```

```
-- the above UPDATE executes after (!) the COMMIT
-- of transaction A
SELECT col FROM t1 WHERE id=1;
-- 102
```

Both UPDATE statements are executed, nothing gets lost.

Please note that transaction B is an example for a 'non-repeatable read' (see below) because the isolation level is '(UN)COMMITTED READ'. First, it reads the value '100' with its SELECT command. Next, it reads '101' with its UPDATE command - after COMMIT of transaction A - and increases it to '102'. If the isolation level would be 'REPEATABLE READ', transaction B would receive the error message 'could not serialize access due to concurrent update' as PostgreSQL's reaction to the UPDATE request.

Non-repeatable read

The example shows a non-repeatable read. (The table *t1* contains its original values.)

Transaction A

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
SELECT col FROM t1 WHERE id=1;
-- 100
```

Transaction B

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
UPDATE t1 SET col=101 WHERE id=1;
SELECT col FROM t1 WHERE id=1;
-- 101
COMMIT;
```

```
SELECT col FROM t1 WHERE id=1;
-- 101 (same transaction, but different value)
-- ' ISOLATION LEVEL REPEATABLE READ' or
-- 'SERIALIZATION' will avoid such an effect
```

Phantom read

The example shows a phantom read. (The table *t1* contains its original values.)

Transaction A

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
SELECT col FROM t1 WHERE id>0;
-- 1 row: 100
```

Transaction B

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
INSERT INTO t1 VALUES (2, 200);
COMMIT;
SELECT col FROM t1 WHERE id>0;
-- 2 rows: 100 and 200
```

```
SELECT col FROM t1 WHERE id>0;
-- 2 rows: 100 and 200
-- (same transaction, same query, but different rows)
-- ' ISOLATION LEVEL SERIALIZABLE'
-- will avoid such an effect
```

Dead lock

The example shows a dead lock. (The table *t1* contains two rows.)

```
DELETE FROM t1;
INSERT INTO t1 VALUES (1, 100);
INSERT INTO t1 VALUES (2, 200);
```

Transaction A

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
UPDATE t1 SET col=col+1 WHERE id=1;
SELECT col FROM t1 WHERE id=1;
-- 101
```

Transaction B

```
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;
UPDATE t1 SET col=col+1 WHERE id=2;
SELECT col FROM t1 WHERE id=2;
-- 201

UPDATE t1 SET col=col+1 WHERE id=1;
.
.
-- must wait for COMMIT/ROLLBACK of transaction A
```

```
UPDATE t1 SET col=col+1 WHERE id=2;
-- must wait for COMMIT/ROLLBACK of transaction B.
--
-- PostgreSQL detects the deadlock and performs a
-- ROLLBACK to overcome the circular situation.
-- message: "ERROR: deadlock detected ..."
```

```
-- processing goes on with a 'success message'  
SELECT col FROM t1 WHERE id>0;  
-- 101  
-- 201  
-- no UPDATES from transaction A. They were  
-- ROLLBACK-ed by PostgreSQL.
```

MVCC

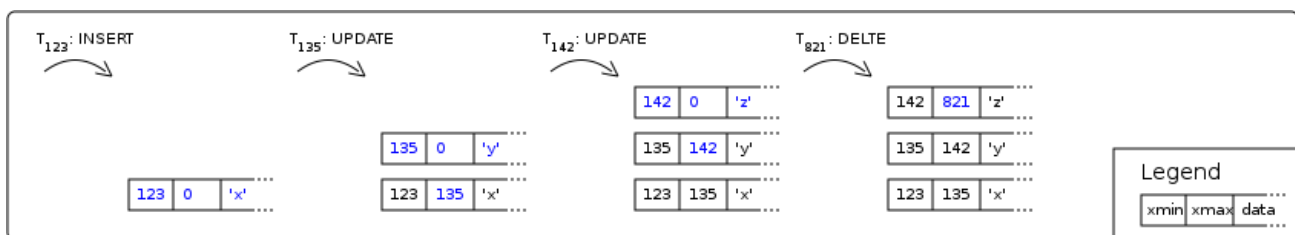
In nearly all cases, PostgreSQL databases must support many clients, which want to add or change data, at the same time. This makes it necessary to protect concurrently running requests from each other - preferably without blocking them. Situations may occur where two clients want to change the same row at the same time or that one client wants to revoke (rollback) his changes while another client may still have tried to read the newest version.

Imagine an Online shop offering the last copy of an article. Two clients display the article at their user interface. After a while, but at the same time, both clients decide to put the article into their shopping cart or even to buy it. Both have seen the article, but only one can be allowed to buy it. The database must enforce an order of the requests, permit the write access to one of them, block the other from writing, and inform the blocked client that the data has been changed by a different process and shall be re-read.

PostgreSQL implements a sophisticated technique to handle concurrent accesses that avoids locking: Multiversion Concurrency Control (MVCC). **Instead of locking a row, the MVCC technique creates a new version of that row when a data change takes place.** "The main advantage of using the MVCC ... rather than locking is that in MVCC locks acquired for querying (reading) data do not conflict with locks acquired for writing data, and so reading never blocks writing and writing never blocks reading. PostgreSQL maintains this guarantee even when providing the strictest level of transaction isolation through the use of ... [Serializable Snapshot Isolation \(SSI\)](#)". ^[1]

The implementation of MVCC is based on transaction IDs (XID). Every transaction in a cluster gets a unique sequential number as its ID. Every INSERT, UPDATE, or DELETE command stores the XID in `xmin` or `xmax` within the affected rows. `xmin`, `xmax`, and some more are system columns contained in every row. Both are not visible with the usual `SELECT * FROM ...` command. But you can read them with commands like `SELECT xmin, xmax, * FROM ...`. The column `xmin` contains the XID of the transaction which has created this version of the row and `xmax` contains the XID of the transaction which has deleted this version, or zero if the version is not deleted.

So, what's going on in detail when write accesses take place? The following graphic shows details concerning `xmin`, `xmax`, and the regular application data.



An INSERT command creates the very first version of a row. Besides its application data 'x', this version contains the ID of the creating transaction 123 in `xmin` and 0 in `xmax`. `xmin` indicates that the version exists since transaction 123 and the value 0 in `xmax` indicates that it is currently not deleted.

Somewhat later, transaction 135 executes an UPDATE of this row by changing the application data from 'x' to 'y'. According to the MVCC principles, the data in the old version of the row is not changed. The value 'x' remains as it was before. Only `xmax` changes to 135. Now, this version is treated as valid exclusively for transactions with XIDs from 123 to 134. In addition to preserve the data in the old version, the UPDATE creates a new version of the complete row with its XID in `xmin`, 0 in `xmax`, and 'y' in the application data (plus all other application data from the old version). This new row version is visible to all future transactions. (Internally, an UPDATE command acts as a DELETE command followed by an INSERT command.)

All subsequent UPDATE commands behave in the same way as the first one: they put their XID in `xmax` of the current version, create a new version with their XID in `xmin` and 0 in `xmax`.

Finally, a row may be deleted by a DELETE command. Even in this case, all versions of the row including the newest one remain in the database - nothing is thrown away. Only `xmax` of the last version is set to the XID of the DELETE transaction, which indicates that it is only visible to transactions with older XIDs - in this example from 142 to 820.

In summary, the MVCC technology creates more and more versions of the same row in the table's heap file and leaves them there, even after a DELETE command. Only the youngest version is relevant for all future transactions. But the system must also preserve some of the older ones for a short time because they could still be requested by transactions that had started before the deleting transaction and hence have a smaller XID. Over time, also the older ones go out of scope for ALL transactions and therefore become ultimately unnecessary. Nevertheless, they do exist physically on the disk and occupy space. They are called *dead rows* and are part of the so-called *bloat*.

Please keep in mind:

- `xmin` and `xmax` indicate the range in which row versions are visible for transactions. This range doesn't imply any direct temporal meaning. The sequence of XIDs reflects only the sequence of transactions' begin events.
- Internally, an UPDATE command acts in the same way as a DELETE command followed by an INSERT command.

- Nothing is removed - with the consequence that the database occupies more and more disk space. It is obvious that this behavior has to be corrected in some way. The next chapter explains how vacuum and autovacuum fulfill this task.

So far this is only a raw description of the principles of MCVV. The implementation considers more problems, e.g.:

- Changes may be revoked by a ROLLBACK command.
- After some time the sequence of XIDs may start from zero (*wrap-around*). In this case xmax can be smaller than xmin.
-

Note

XIDs are sequences (with a reserved value to handle *wrap-around* in pre-9.4 PostgreSQL versions). PostgreSQL knows some configuration parameters concerning transactions and their XIDs with names like *xxx_age*, e.g.: *vacuum_freeze_min_age*. For such parameters, the 'age' doesn't specify a period of time but represents a certain number of transactions, e.g., 100 millions.

References

1. MVCC in PostgreSQL [<https://www.postgresql.org/docs/current/mvcc-intro.html>]

Vacuum

Eliminating Bloat

As we have seen in the [MVCC](#) chapter, the database tends to occupy more and more disk space caused by *bloat*: over time more and more logically deleted but physically existing old versions of rows arise within heap and index files. This chapter explains how the SQL command **VACUUM** and the automatically running **Autovacuum processes** clean up files and thereby prevent their endless growth.

One process of the Instance is the Autovacuum daemon. It continuously monitors the state of all databases based on values that are collected by the Statistics Collector and starts Autovacuum processes whenever it detects certain situations, e.g.: a huge number of modifications within a table. This leads to the intended dynamic behavior of PostgreSQL: Only when it is necessary, Autovacuum cleans up the files. In addition, client processes can issue the SQL command **VACUUM** at any time. DBAs do this interactively when they recognize critical situations, or they start it in periodically running batch jobs. In most cases, this is not necessary because of the constantly running Autovacuum daemon.

The central value to determine which of the physically existing row versions are no longer needed is `xmax`, which shows what transaction has deleted the row. The elimination operation must evaluate it against several criteria which must all apply:

- `xmax` must be different from zero because a value of zero indicates that the row version is not deleted.
- `xmax` must contain an `XID` which is older than the oldest `XID` of all currently running transactions. That guarantees that no existing or upcoming transaction will have read or write access to this row version.
- The transaction of `xmax` must be committed. If it is still running or was rollback-ed, this row version is treated as valid (not deleted).
- If there is a situation that the row version is part of multiple transactions, more actions must be taken.

When the vacuum operation detects such an outdated row version, it marks its space as free for future use of write actions, optimizes the physical arrangement of the remaining versions on the page, and removes index tuples pointing to the removed row version. But only in rare situations (or in the case of **VACUUM FULL**), this space is released to the operating system. In most cases, it remains occupied by the database and will be used by future **INSERT**, **UPDATE** or **DELETE** commands. Hence, even with successful running Autovacuum, the size of files does not shrink; but they have more respectively huger 'holes' where coming data can be stored. Only after this 'hole-space' is exhausted (per page), it gets necessary to claim new disc space from the operating system to store new data.

An exception to this conservative behavior is the SQL command **VACUUM FULL**. It creates a new file at the operating system level, copies all valid row versions to it with no extra space by ignoring the 'holes', and deletes the old file. But it is slower and requires an exclusive lock on the affected tables.

Because vacuum operations typically are *I/O* intensive, which can hinder other activities, Autovacuum avoids performing many vacuum operations in bulk. Instead, it carries out many small actions with time delays in between. The SQL command **VACUUM** runs immediately and without any time delay.

More Actions

VACUUM, as well as Autovacuum, don't just eliminate *bloat*. They perform additional tasks for minimizing future *I/O* activities of themselves as well as of other processes. This extra work can be done in a very efficient way since in most cases the expensive physical access to pages has taken place anyway. The additional operations are:

- *Freeze*: It marks certain row versions as frozen. This means that they are treated as 'valid' (visible) forever, independent from the wraparound problem (see later).
- *Visibility Map and Free Space Map*: It logs information about the state of the handled pages in two additional files, the *Visibility Map* and the *Free Space Map*.
- *Statistics*: Similar to the *Statistics Collector* it collects statistics about the number of rows per table, the distribution of values, and so on, as the basis for decisions of the query planner.

External links

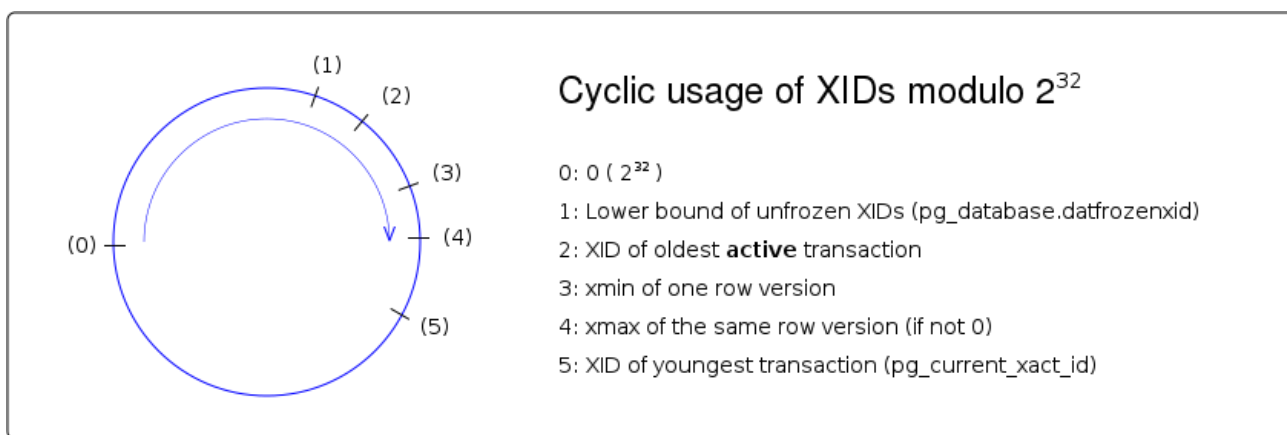
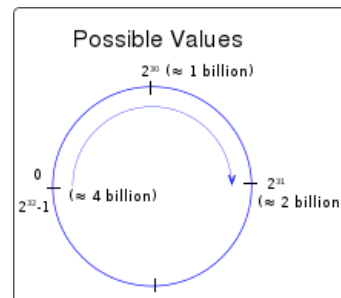
PostgreSQL Documentation concerning **VACUUM** (<https://www.postgresql.org/docs/current/routine-vacuuming.html>)

Wraparound and Freeze

Two Fundamental Problems

Transactions are identified by ids which are realized as unsigned 32-bit integers and named XID. Transactions and their XIDs are known at the cluster level, covering all databases. Similar to sequences, XIDs are incremented by +1 for every new transaction. Sooner or later, this limited space of 2^{32} numbers is exhausted, and it becomes necessary to restart the sequence from the beginning (the values 0, 1, and 2 are skipped because they are reserved for particular purposes). This restart of XIDs is called a *wraparound* and each cycle an *epoch*.

It's unlikely that more than 2^{32} transactions exist in a cluster at the same time or that a single transaction lasts for such a long time that its XID collides with the same value of the next cycle. At first glance, this cyclic usage of the ' 2^{32} Universe' seems to be safe and easy to implement. Nevertheless, huge problems will arise with this simple strategy. The reason is that XIDs are stored in system columns within every row (see `xmin`, `xmax` in the [MVCC](#) chapter). And rows stay for a very long time in the database, in many cases forever.



XID Collision

The first problem is that after a wraparound, the next XIDs (3, 4, 5, ...) may collide with XIDs of the previous epoch. They are no longer unique because system columns of old rows may contain the same values. But transactions must be able to decide whether retrieved rows are modified (by other transactions) after their own start-time or a long time ago. We call this first problem **XID Collision**.

Sudden Death

The second problem correlates with [MVCC](#) and the timeline of transactions. Rows may exist in multiple versions. When a transaction modifies a row and stays alive a little longer - no COMMIT or ROLLBACK because of more activities -, other processes shall 'see' the version of the row as of the start-time of the transaction and not the uncommitted modification. Hence, a mechanism must hide the ongoing changes and give other transactions the feeling of a stable data environment. The system realizes this by considering additional criteria with every SQL command, especially - but not only - the system column `xmin`.

You can imagine of those additional criteria that the system silently supplements every query with the predicate `xmin < my_xid`. (This is only an illustration in pseudo-code, the real implementation is different.) It guarantees that changes happening after the start of the requesting transaction are invisible to it.

So far, so good.

But what will happen after a wraparound? The next transactions will have very small XIDs, e.g., '5'. And what will be the result of an `xmin < 5`? Near nothing. All rows with `xmin` between 5 and $2^{32} - 1$ are no longer part of any query. In contrast to the situation a few moments ago, all data suddenly disappeared. It is still in the database, but it is unreachable. We call this second problem **Sudden Death**.

Solution

Step 1: At the conceptual level, the full 2^{32} Universe' is divided into two halves of 2^{31} numbers. One split point is the current transaction id `pg_current_xact_id` (it was called `txid_current` before PostgreSQL version 13) and the other is the opposite side of the circle `pg_current_xact_id + 231` (or `pg_current_xact_id - 231` what is the same). So, the split points are not fixed values but follow dynamically the ongoing of new transactions. One halve represents the previously used and therefore exhausted XIDs; the other halve such XIDs, which are - per definition - free. They will be allocated in the future. Please note the dynamic aspect: with every new transaction in the cluster `pg_current_xact_id` and the border between 'past / future' moves forward. This is a metaphor of an endless walk through time where after some time the old problems will be forgotten or at least idealized.

The idea can be realized by a modification of the above `xmin < my_xid` predicate to an `if/else` block:

```
if (my_xid < 231)
  return rows with: xmin < my_xid OR xmin > my_xid + 231
else
  return rows with: xmin < my_xid AND xmin > my_xid + 231
```

Of course, this is a simplification and many other criteria like COMMIT status, 'is deleted', and other things must be considered. It focuses purely on the aspect of the 'past / future' metaphor.

Note: With this algorithm the 'critical point' changes from 0 resp. 2^{32} to `pg_current_xact_id + 231`. It is called the *wraparound point* and the line between `pg_current_xact_id + 231` and `pg_current_xact_id` the *wraparound horizon*.

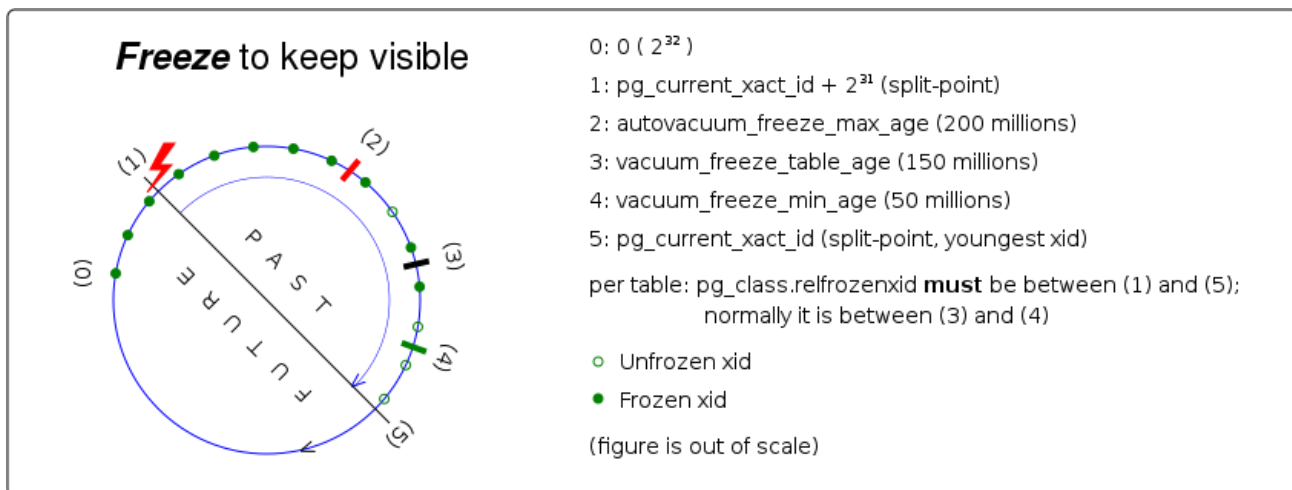
Step 2: The outlined algorithm ensures the visibility of 50% of all possible XIDs. But what's going on with the others? As mentioned, rows may stay in the database forever holding very old XIDs in `xmin`. This halve must also be considered. The idea of how to access the complete range of possible XIDs is to complement the previous algorithm with the introduction of a flag that marks certain rows as 'visible forever' (respectively visible until the next write operation to them). This marking is not possible as long as one or more transactions potentially get write access to them. Fortunately, the sequence of new XIDs goes strictly forward, and overtime transactions with old XIDs finish. PostgreSQL does not only know, which is the current XID `pg_current_xact_id`, but also which is the **oldest active** XID per connection (`pg_stat_activity.backend_xmin`), and which is the lower bound of all unfrozen XIDs per table (`pg_class.relfrozensxid`) and per database (`pg_database.datfrozensxid`). Rows with `xmin` older than `oldest` (`pg_stat_activity.backend_xmin`) are candidates for such a flag. No running transaction has or will get write access to them, only newer ones. According to MVCC, they will create a new row version, this one keeps as it is.

It is one of the two main duties of VACUUM to perform this freezing. It marks the identified rows with a flag in their header `t_infomask` as 'visible forever'. From this point on no comparisons with `xmin` take place. The rows are always treated as visible, even if they are part of the 'future'. This marking is called FREEZE and the status of the row FROZEN.

Now the algorithm for retrieving rows changes to:

```
-- 'frozen' rows will always be returned
if (my_xid < 231)
  return rows with: frozen OR
                  (xmin < my_xid OR xmin > my_xid + 231)
else
  return rows with: frozen OR
                  (xmin < my_xid AND xmin > my_xid + 231)
```

With this extension, the two mentioned problems are solved. The system can generate XIDs even after a wraparound without risk of collision with old XIDs. The old ones may exist, but they are not touched in any way. Second, the algorithm finds all relevant XIDs whether there was a wraparound or not.



Wraparound Failure

It is possible that a transaction - intentionally or by an error in an application - stays alive for a long time. Over time its *XID* becomes the oldest one in the complete cluster and can be retrieved from `pg_stat_activity.backend_xmin`. As long as this situation continues, the gap between the ongoing *wrapping point* `pg_current_xact_id + 231` and `pg_stat_activity.backend_xmin` gets smaller and smaller. If the gap would close completely, we would see all the problems described at the beginning of the chapter. This is called *wraround failure* and must be avoided under all circumstances. VACUUM is doing its best to freeze as many rows as possible. But if a long-living transaction prevents freezing and the size of the gap falls below a certain limit, VACUUM runs in an 'aggressive mode' and works on all pages of affected tables, independent from the above-mentioned values; if also this fails, the cluster stops the creation of new transactions and prevents further write actions.

Details

Note: These details can be skipped by a novice reader without losing the context of the ongoing chapters.

To freeze any row version, VACUUM must check several criteria:

- `xmax` must be zero because only non-deleted rows can be visible forever.
- `xmin` must be older than all currently existing transactions `oldest(pg_stat_activity.backend_xmin)`. This guarantees that no existing transaction can modify or delete the version.
- The transactions of `xmin` and `xmax` must be committed.

At what point in time does the freeze operation take place? Please note that there are configuration parameters with names like `xxx_age`. They define distances - mostly to `pg_current_xact_id` -, where the actions of VACUUM shall start. 'age' in this context doesn't imply a certain period of time, it's always a pure number that counts transactions, e.g., 50 million. Please also note, that VACUUM always reads complete physical pages and works on the row versions found there.

- When a client issues the SQL command VACUUM with its FREEZE option. In this case, all pages of the affected tables are processed that are marked in the [Visibility Map](#) as potentially having unfrozen rows.
- When a client issues the SQL command VACUUM without any option and there are XIDs older than `vacuum_freeze_table_age` (default: 150 million) minus `vacuum_freeze_min_age` (default: 50 million). As before, all pages are processed that are marked in the [Visibility Map](#) to potentially have unfrozen rows.
- When an Autovacuum process runs. Such a process acts in one of two modes: In the *normal mode* it skips pages with row versions that are younger than `vacuum_freeze_min_age` (default: 50 million) and works only on pages where all XIDs are older. The skipping of young XIDs prevents work on such pages, which are likely to be changed by one of the future SQL commands. The process switches to an *aggressive mode* if it recognizes that for the processed table the oldest XID exceeds `vacuum_freeze_table_age` (default: 150 million). In this *aggressive mode*, Autovacuum processes all pages of the affected table.

VACUUM and Autovacuum know to which value the oldest unfrozen XID has moved forward per table and logs the value in `pg_class.relFrozenxid`. The distance between this value and the `pg_current_xact_id` split point becomes smaller (there are potentially unfrozen rows), and the distance to the *wraround point* `pg_current_xact_id + 231` becomes larger (there are only frozen rows). That is how the freezing follows the moving 'past' / 'future' horizon.

Note: Before version 9.4 of PostgreSQL, the freeze algorithm stored the value '2' (FrozenTransactionId) in `xmin` instead of setting a flag in `t_infomask`.

Visibility Map and Free Space Map

Every table is stored in a separate disk file. The names of such files consist of numbers which are the internal Object Identifiers (OID) of the table used in the System Catalog. Each such file is accompanied by a file for its *Visibility Map* and another one for its *Free Space Map*. They have the same names expanded by the suffix '_vm' respectively '_fsm'. An example for such a triplet of filenames is: 3083, 3083_vm, and 3083_fsm.

The two additional files contain metainformation to optimize I/O activities to the original file - especially for vacuuming and freezing, but also for other write activities. Because I/O works per physical page, the metainformation addresses complete physical pages, not any part of pages like rows or versions.

Visibility Map

The Visibility Map contains two flags — stored as two bits — for each page of the original file. The first bit indicates that the associated page contains only valid row versions, i.e. there is no bloat to be vacuumed. The second bit indicates that the page only contains already frozen row versions.

Please consider two details. First, in most cases a page contains many rows or row versions. However, both flags are associated with the page, not with an individual row or row version. The flags are set only under the condition that they are valid for ALL row versions stored on the page. Second, since there are only two bits per page (2 bits correlate to 8 kiloBytes, which is a relation of about 1 : 32.000), the Visibility Map is considerably smaller than the original file.

VACUUM and Autovacuum set the flags. Every write operation on any row version of the page clears the flags.

The Visibility Map helps VACUUM and Autovacuum to save unnecessary I/O. When the first bit is set, it's unnecessary to read the original page and check its content for removing bloat. It's clear that there is no bloat. Correspondingly, if VACUUM or Autovacuum have to perform freezing of rows, they can skip pages where the second bit signals that the page only contains already frozen row versions - no freeze is necessary.

Free Space Map

The Free Space Map tracks the amount of free, unused space per page. It is organized as a highly condensed B-tree of (rounded) free space size per page.

VACUUM and Autovacuum change the Free Space Map according to their write operations (marking of row versions as 'obsolete' and rearranging the physical layout of pages). Other write operations consult the Free Space Map to locate pages with enough free space for the intended write operations and change the Free Space Map afterward.

WAL

Usage

WAL (Write Ahead Logging) files are files, where changed data values are stored in a binary format. This is additional information and in this respect it is redundant to the information in the database files. WAL files are a specific kind of 'diff' files.

Writing to WAL files is very fast as they are written always sequentially. In contrast to WAL files database files are organized in special structures like trees, which possibly must be reorganized during write operations or which points to blocks at far positions. Thus writes to database files are much slower.

For the mentioned performance reasons, when a client requests a write operation like `UPDATE` or `DELETE` the modifications to the data are done in a special sequence and - in some parts - asynchronously to the client requests. First, data is written and flushed to WAL files. Second, it is stored in shared buffers in RAM. Finally, it is written from shared buffers to database files. The client doesn't wait until the end of all operations. After the first two very fast actions, he is informed that his request is completed. The third operation is performed asynchronously at an later (or prior) point in time.

Remove

WAL files are collected in the directory `pg_wal` (`pg_xlog` in PostgreSQL versions prior to version 10). Depending on the write activities on the database the total size of all WAL files may increase dramatically. Therefore the system must delete them when they are no longer needed. WAL files are available for deletion after the changes in the shared buffers (which correlate to the content of the WAL files) are flushed to the database files. As it is guaranteed that this criterion is met after a `CHECKPOINT`, there are some dependencies between WAL file delete operations and `CHECKPOINTS`:

- You can define a limit for the total size of all files in the directory: `max_wal_size`. If it is reached, PostgreSQL performs an automatic `CHECKPOINT` operation.
- You can define a `checkpoint_timeout` in seconds. No later than this number of seconds, PostgreSQL performs an automatic `CHECKPOINT` operation.

In both cases the shared buffers gets written to disc, a checkpoint-record is written to the actual WAL file and all older WAL files are ready to be deleted.

The deletion of WAL files may be prevented by other criterion, especially by failing archive commands. `max_wal_size` is a soft limit and may be exceeded in such situations.

ClientServerComm Client

Before a client program like `createdb`, `psql`, `pg_dump`, `vacuumdb`, ... can perform any action on a database, it must establish a connection to that database (or cluster). To do so, it must provide concrete values for the essential boundary conditions.

- The IP address or DNS name of the server, where the instance is running.
- The port on this server, to whom the instance is listening.
- The name of the database within the cluster (= IP/port combination).
- The name of the user (= role) with which the client program wants to work
- The password of this user.

You can specify these values in three different ways:

- as explicit parameters of the client program
- as environment variables
- as a fixed line of text in the special file `pgpass`.

Parameters

You can specify the parameters in the usual short (-) or long (--) format of `createdb`, `psql`, `pg_dump`, `vacuumdb`, and other standard PostgreSQL command line tools.

```
$ # Example
$ psql -h www.dbserver.com --port=5432 ....
```

The parameter names and their meanings are:

Short Form	Long Form	Meaning
-h	--host	IP or DNS
-p	--port	port number (default: 5432)
-d	--dbname	database within the cluster
-U	--username	name of the user

If necessary, the client program will prompt for the password.

Environment Variables

As an alternative to the parameter passing you can define environment variables within your shell.

Environment Variable	Meaning
PGHOST	IP or DNS
PGPORT	port number (default: 5432)
PGDATABASE	database within the cluster
PGUSER	name of the user
PGPASSWORD	password of this user (not recommended)
PGPASSFILE	name of a file where those values are stored as plain text, see below (default: <code>.pgpass</code>)

File 'pgpass'

Instead of using parameters or environment variables as shown above you can store those values in a file. Use one line per definition in the form:

```
hostname:port:database:username:password
```

The default filename on UNIX systems is `~/.pgpass` and on Windows: `C:\Users\MyUser\AppData\Roaming\postgresql\pgpass.conf`. On UNIX systems the file protections must disallow any access of world or group: `chmod 0600 ~/.pgpass`.

You can create the file with any text editor. This is not necessary if you use pgAdmin. pgAdmin creates the file automatically after a successful connection and stores the actual connection values.

ClientServerComm

Protocol

All access to data is done by server (or backend) processes, to which client (or frontend) processes must connect to. In most cases instances of the two process classes reside on different hardware, but it's also possible that they run on the same computer. The communication between them uses a PostgreSQL-specific protocol, which runs over TCP/IP or over UNIX sockets. It is implemented in the C library *libpq*. For every incoming new connection the backend process (sometimes called the *postmaster*- process) creates a new *postgres* backend process. This backend process gets part of the *PostgreSQL instance*, which is responsible for data accesses and database consistency.

The protocol handles the authentication process, client request, server responses, exceptions, special situations like a NOTIFY, and the final regular or irregular termination of the connection.

Driver

Most client programs - like *psql* - use this protocol directly. Drivers like ODBC, JDBC (type 4), Perl DBI, and those for Python, C, C++, and much more are also based on *libpq*.

You can find an extensive list of drivers at the postgres wiki [\[1\]](#) and some more commercial and open source implementations at the 'products' site [\[2\]](#).

Authentication

Clients must authenticate themselves before they get access to any data. This process has one or two stages. During the first - optional - step the client gets access to the server by satisfying the operating system hurdles. This is often realized by delivering a public ssh key. The authentication with PostgreSQL is a separate, independent step using a database-username, which may or may not correlate to an operating system username. PostgreSQL stores all rules for this second step in the file *pg_hba.conf*.

pg_hba.conf stores every rule in one line, one rule per line. The lines are evaluated from the top of *pg_hba.conf* to bottom and the first matching line applies. The main layout of these lines is as follows

```
-----  
local DATABASE USER METHOD [OPTIONS]  
host DATABASE USER ADDRESS METHOD [OPTIONS]  
-----
```

Words in upper case must be replaced by specific values. Lower case words like *local* and *host* are key words. They decide, for which kind of connection the rule shall apply: *local* for clients residing at the same computer as the backend (they use UNIX sockets for the communication) and *host* for clients at different computers (they use TCP/IP). There is one notable exception. In the former case clients can use the usual TCP/IP syntax `--host=localhost --port=5432` to switch over to use TCP/IP. Thus the *host* syntax applies for them.

DATABASE and USER have to be replaced by the name of the database and the name of the database-user, for which the rule will apply. In both cases the key word ALL is possible to define, that the rule shall apply to all databases and respectively all database-users.

ADDRESS must be replaced by the hostname or the IP adress plus CIDR mask of the client, for which the rule will apply. IPv6 notation is supported.

METHOD is one of the following. The thereby defined rule (=line) applies, if database/user/address combination is the first matching combination in *pg_hba.conf*.

- trust: The connection is allowed without a password.
- reject: The connection is rejected.
- password: The client must send a valid user/password combination.
- md5: Same as 'password', but the password is encrypted.
- ldap: It uses LDAP as the password verification method.
- peer: The connection is allowed, if the client is authorized against the operation system with the same username as the given database username. This method is only supported on local connections.

There are some more techniques in respect to the METHOD.

Some examples:

```
-----  
# joe cannot connect to mydb - eg. with psql -, when he is logged in to the backend.  
local mydb joe reject  
-----
```

```
# bill (and all other persons) can connect to mydb when they are logged in to the
# backend without specifying any further password. joe will never reach this rule, he
# is rejected by the rule in the line before. The rule sequence is important!
local mydb all trust

# joe can connect to mydb from his workstation '192.168.178.10', if he sends
# the valid md5 encrypted password
host mydb joe 192.168.178.10/32 md5

# every connection to mydb coming from the IP range 192.168.178.0 - 192.168.178.255
# is accepted, if they send the valid md5 encrypted password
host mydb all 192.168.178.0/24 md5
```

For the DATABASE specification there is the special keyword REPLICATION. It denotes the streaming replication process. REPLICATION is not part of ALL and must be specified separately.

References

1. Driver Wiki [10] (https://wiki.postgresql.org/wiki/List_of_drivers)
2. Commercial and open source driver [11] (<http://www.postgresql.org/download/products/2/>)

Security

Roles

PostgreSQL supports the concept of roles ^[1] to handle security issues within the database. Roles are independent from operating system user accounts (with the exception of the special case peer authentication which is defined in the `pg_hba.conf` file).

The concept of roles subsumes the concepts of individual users and groups of users with similar rights. A role can be thought of as either an individual database user, or a group of database users, depending on how the role is set up. Thus the outdated SQL command `CREATE USER . . .` is only an alias for `CREATE ROLE . . .`. Roles have certain privileges on database objects like tables or functions and can assign those privileges to other roles. Roles are global across a cluster - not per individual database.

Often individual users, which shall have identical privileges, are grouped together to a user group and the privileges are granted to that group.

```
-- ROLE, in the sense of a group of individual users or other roles
CREATE ROLE group_1 ENCRYPTED PASSWORD 'xyz';
-- assign some rights to the role
GRANT SELECT ON table_1 TO group_1;
-- ROLE, in the sense of some individual users
CREATE ROLE adam LOGIN ENCRYPTED PASSWORD 'xyz'; -- Default is NOLOGIN
CREATE ROLE anne LOGIN ENCRYPTED PASSWORD 'xyz';
-- the link between user group and individual users
GRANT group_1 TO adam, anne;
```

With the `CREATE ROLE` command you can assign the privileges `SUPERUSER`, `CREATEDB`, `CREATEROLE`, `REPLICATION` and `LOGIN` to that role. With the `GRANT` command you can assign access privileges to database objects like tables. The second purpose of the `GRANT` command is the definition of the group membership.

In addition to the roles created by the database administrator there is always the special role `PUBLIC`, which can be thought of as a role which is a member of all other roles. Thus, privileges assigned to `PUBLIC` are implicitly given to all roles, even if those roles are created at a later stage.

List existing roles

Roles can be listed with the following commands.

With SQL, this will display an additional set of postgresQL default roles that group together sets of common access levels:

```
SELECT rolname FROM pg_roles;
```

or the `psql` command:

```
\du
```

Users

```
select * from postgres.pg_catalog.pg_user
```

References

1. Concept of roles ^[12] (<http://www.postgresql.org/docs/current/static/user-manag.html>)

Replication

Replication is the process of transferring data changes from one or many databases (master) to one or many other databases (standby) running on one or many other nodes. The purpose of replication is

- High Availability: If one node fails, another node replaces him and applications can work continuously.
- Scaling: The workload demand may be too high for one single node. Therefore, it is spread over several nodes.

Concepts

PostgreSQL offers a bunch of largely mutually independent concepts for use in replication solutions. They can be picked up and combined - with only few restrictions - depending on the use case.

Events

- With *Trigger Based Replication* a trigger (per table) starts the transfer of changed data. This technique is outdated and not used.
- With *Log Based Replication* such information is transferred, which describes data changes and is created and stored in WAL files anyway.

Shipping

- *WAL-File-Shipping Replication* (or *File-based Replication*) denotes the transfer of completely filled WAL files (16 MB) from master to standby. This technique is not very elegant and will be replaced by *Streaming Replication* over time.
- *Streaming Replication* denotes the transfer of log records (single change information) from master to standby over a TCP connection.

Primary parameter: 'primary_conninfo' in recovery.conf on standby server.

Format

- In *Physical Format* the transferred WAL records have the same structure as they are used in WAL files. They reflect the structure of database files including block numbers, VACUUM information and more.
- The *Logical Format* is a decoding of WAL records into an abstract format, which is independent from PostgreSQL versions and hardware platforms.

Primary parameter: 'wal_level=logical' in postgres.conf on master server.

Synchronism

- In *Asynchronous Replication* data is transferred to a different node without waiting for a confirmation of its receiving.
- In *Synchronous Replication* the data transfer waits - in the case of a COMMIT - for a confirmation of its successful processing on the standby.

Primary parameter: 'synchronous_standby_names' in postgres.conf on master server.

Standby Mode

- *Hot*: In *Hot Standby Mode* the standby server runs in 'recovery mode', accepts client connections, and processes their read-only queries.
- *Warm*: In *Warm Standby Mode* the standby server runs in 'recovery mode' and doesn't allow clients to connect.
- *Cold*: Although it is not an official PostgreSQL term, *Cold Standby Mode* can be associated with a not running standby server with log-shipping technique. The WAL files are transferred to the standby but not processed until the standby starts up.

Primary parameter: 'hot_standby=on/off' in recovery.conf on standby server.

Architecture

In contrast to the above categories, the two different architectures (*Master/Standby* and *Multi-Master*) are not strictly distinct from each other. For example, if you focus on atomic replication channels of a *Multi-Master* architecture, you will also see a *Master/Standby* replication.

- The *Master/Standby* architecture denotes a situation, where one or many standby nodes receive change data from one master node. In such situations standby nodes may replicate the received data to other nodes, so they are master and standby at the same time.
- The *Multi-Master* architecture denotes a situation, where one or many standby nodes receive change data from many master nodes.

Configuration

There are 3 main configuration files:

- 'postgres.conf'
- 'pg_hba.conf'
- 'recovery.conf'

The 'postgres.conf' is the main configuration file. It is used to configure the master site. Additional instances can exist on standby sites.

The 'pg_hba.conf' is the security and authentication configuration file.

The 'recovery.conf' was optional and contained restore and recovery configurations. As of PostgreSQL-12 it is no longer used and its existence will prevent the server from starting. The recovery.conf settings as of PostgreSQL-12 can be set in postgres.conf.

Because the great number of possible combinations of concepts and correlating configuration values may be confusing at the beginning, this book will focus on a minimal set of initial configuration values.

Shipping: WAL-File-Shipping vs. Streaming

WAL files are generated anyway because they are necessary for recovery after a crash. If they are - additionally - used to ship information to a standby server, it is necessary to add some more information to the files. This is activated by choosing 'replica' or 'logical' as a value for wal_level.

```
# WAL parameters on MASTER's postgres.conf
wal_level=replica          # 'archive' | 'hot_standby' in versions prior to PG 9.6
archive_mode=on           # activate the feature
archive_command='scp ...' # the transfer-to-standby command (or to an archive location, which is the original
                           # purpose of this command)
```

If you switch the shipping technique to streaming instead of WAL-file you must not deactivate WAL-file generating and transferring. For safety reasons you may want to transfer WAL files anyway (to a platform different from the standby server). Therefore, you can retain the above parameters in addition to streaming replication parameters.

The streaming activities are initiated by the standby server. When he finds the file 'recovery.conf' during its start up, he assumes that it is necessary to perform a recovery. In our case of replication he uses nearly the same techniques as in the recovery-from-crash situation. The parameters in 'recovery.conf' advice him to start a so-called WAL receiver process within its instance. This process connects to the master server and initiates a WAL sender process over there. Both exchange information in an endless loop whereas the standby server keeps in 'recovery mode'.

The authorization at the operating system level shall be done by exchanging ssh keys.

```
# Parameters in the STANDBY's recovery.conf
standby_mode=on # activates standby mode
# How to reach the master:
primary_conninfo='user=<replication_dbuser_at_master> host=<IP_of_master_server> port=<port_of_master_server>
                 sslmode=prefer sslcompression=1 krbsrvname=...'
# This file can be created by the pg_basebackup utility, see below
```

On the master site there must be a privileged database user with the special role REPLICATION:

```
CREATE ROLE <replication_dbuser_at_master> REPLICATION ...;
```

And the master must accept connections from the standby in general and with a certain number of processes.

```
# Allow connections from standby to master in MASTER's postgres.conf
listen_addresses = '<ip_of_standby_server>' # what IP address(es) to listen on
max_wal_senders = 5 # no more replication processes/connections than this number
```

Additionally, authentication of the replication database user must be possible. Please notice that the key word ALL for the database name does not include the authentication of the replication activities. 'Replication' is a key word of its own and must be noted explicitly.

```
# One additional line in MASTER's pg_hba.conf
# Allow the <replication_dbuser> to connect from standby to master
host replication <replication_dbuser> <IP_of_standby_server>/32 trust
```

Now you are ready to start. First, you must start the master. Second, you must transfer the complete databases from the master to the standby. And at last you can start the standby. Just as the replication, the transfer of the databases is initiated at the standby site.

```
pg_basebackup -h <IP_of_master_server> -D main --wal-methode=stream --checkpoint=fast -R
```

The utility `pg_basebackup` transfers everything to the directory 'main' (shall be empty), in this case it uses the streaming method, it initiates a checkpoint at the master site to enforce consistency of database files and WAL files, and due to the `-R` flag it generates previously mentioned `recovery.conf` file.

Format: Physical vs. Logical

The decoding of WAL records from their physical format to a logical format was introduced in PostgreSQL 9.4. The physical format contains - among others - block numbers, VACUUM information and it depends on the used character encoding of the databases. In contrast, the logical format is independent from all these details - conceptually even from the PostgreSQL version. Decoded records are offered to registered streams for consuming.

This logical format offers some great advantages: transfer to databases at different major release levels, at different hardware architectures, and even to other writing master. Thus multi-master-architectures are possible. And additionally it's not necessary to replicate the complete cluster: you can pick single database objects.

In release 9.5 the feature is not delivered with core PostgreSQL. You must install some extensions:

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

As the feature is relative new, we don't offer details and refer to the [documentation \(http://www.postgresql.org/docs/current/static/logicaldecoding.html\)](http://www.postgresql.org/docs/current/static/logicaldecoding.html). And there is an important project [Bi-Directional Replication \(http://bdr-project.org/docs/stable/index.html\)](http://bdr-project.org/docs/stable/index.html), which is based on this technique.

Synchronism: synchron vs. asynchron

The default behaviour is asynchronous replication. This means that transferred data is processed at the standby server without any synchronization with the master, even in the case of a COMMIT. In opposite to this behaviour the master of a synchronous replication waits for a successful processing of COMMIT statements at the standby before he confirms it to its client.

The synchronous replication is activated by the parameter 'synchronous_standby_names'. Its values identify such standby servers, for which the synchronicity shall take place. A '*' indicates all standby server.

```
# master's postgres.conf file  
synchronous_standby_names = '*'
```

Standby Mode: hot vs. warm

As long as the standby server is running, he will continuously handle incoming change information and store it in its databases. If there is no necessity to process requests from applications, he shall run in warm standby mode. This behaviour is enforced in the `recovery.conf` file.

```
# recovery.conf on standby server  
hot_standby = off
```

If he shall allow client connections, he must start in hot standby mode. In this mode read-only access from clients are possible - write actions are denied.

```
# recovery.conf on standby server  
hot_standby = on
```

To generate enough information on the master site for the standby's hot standby mode, its WAL level must also be replica or higher.

```
# postgres.conf on master server  
wal_level = replica
```

Typical Use Cases

We offer some typical combinations of the above-mentioned concepts and show its advantages and disadvantages.

Warm Standby with Log-Shipping

In this situation a master sends information about changed data to a standby using completely filled WAL files (16 MB). The standby continuously processes the incoming information, which means that the changes made on the master are seen at the standby over time.

To build this scenario, you must perform steps, which are very similar to [Backup with PITR](#):

- Take a physical backup exactly as described in [Backup with PITR](#) and transfer it to the standby.
- At the master site `postgres.conf` must specify `wal_level=replica`; `archive_mode=on` and a copy command to transfer WAL files to the standby site.
- At the standby site the central step is the creation of a `recovery.conf` file with the line `standby_mode='on'`. This is a sign to the standby to perform an 'endless recovery process' after its start.
- `recovery.conf` must contain some more definitions: `restore_command`, `archive_cleanup_command`

With this parametrisation the master will copy its completely filled WAL files to the standby. The standby processes the received WAL files by copying the change information into its database files. This behaviour is nearly the same as a recovery after a crash. The difference is, that the recovery mode is not finish after processing the last WAL file, the standby waits for the arrival of the next WAL file.

You can copy the arising WAL files to a lot of servers and activate warm standby on each of them. Doing so, you get a lot of standbys.

Hot Standby with Log-Shipping

This variant offers a very valuable feature in comparison with the warm standby scenario: applications can connect to the standby and send read requests to him while he runs in standby mode.

To achieve this situation, you must increase `wal_level` to `hot_standby` at the master site. This leads to some additional information in the WAL files. And on the standby site you must add `hot_standby=on` in `postgres.conf`. After its start the standby will not only process the WAL files but also accept and response to read-requests from clients.

The main use case for hot standby is load-balancing. If there is a huge number of read-requests, you can reduce the masters load by delegating them to one or more standby servers. This solution scales very good across a great number of parallel working standby servers.

Both scenarios *cold/hot with log-shipping* have a common shortage: The amount of transferred data is always 16 MB. Depending on the frequency of changes at the master site it can take a long time until the transfer is started. The next chapter shows a technique which does not have this deficiency.

Hot Standby with Streaming Replication

The use of files to transfer information from one server to another - as it is shown in the above log-shipping scenarios - has a lot of shortages and is therefore a little outdated. Direct communication between programs running on different nodes is more complex but offers significant advantages: the speed of communication is incredible higher and in much cases the size of transferred data is smaller. In order to gain these benefits, PostgreSQL has implemented the [streaming replication technique](#), which connects master and standby servers via TCP. This technique adds two additional processes: the *WAL sender* process at the master site and the *WAL receiver* process at the standby site. They exchange information about data changes in the master's database.

The communication is initiated by the standby site and must run with a database user with REPLICATION privileges. This user must be created at the master site and authorized in the master's `pg_hba.conf` file. The master must accept connections from the standby in general and with a certain number of processes. The authorization at the operating system level shall be done by exchanging ssh keys.

```
Master site:
=====

-- SQL
CREATE ROLE <replication_dbuser_at_master> REPLICATION ...;

# postgresql.conf: allow connections from standby to master
listen_addresses = '<ip_of_standby_server>' # what IP address(es) to listen on
max_wal_senders = 5 # no more replication processes/connections than this number
# make hot standby possible
wal_level = replica # 'hot_standby' in versions prior to PG 9.6

# pg_hba.conf: one additional line (the 'all' entry doesn't apply to replication)
# Allow the <replication_dbuser> to connect from standby to master
host replication <replication_dbuser> <IP_of_standby_server>/32 trust

Standby site:
=====

# recovery.conf (this file can be created by the pg_basebackup utility, see below)
standby_mode=on # activates standby mode
# How to reach the master:
primary_conninfo='user=<replication_dbuser_at_master_server> host=<IP_of_master_server> port=<port_of_master_server>
sslmode=prefer sslcompression=1 krbsrvname=...'

# postgres.conf: activate hot standby
hot_standby = on
```

Now you are ready to start. First, you must start the master. Second, you must transfer the complete databases from the master to the standby. And at last you start the standby. Just as the replication activities, the transfer of the databases is initiated at the standby site.

```
pg_basebackup -h <IP_of_master_server> -D main --wal-method=stream --checkpoint=fast -R
```

The utility *pg_basebackup* transfers everything to the directory 'main' (shall be empty), in this case it uses the streaming method, it initiates a checkpoint at the master site to enforce consistency of database files and WAL files, and due to the -R flag it generates the previous mentioned recovery.conf file.

The activation of the 'hot' standby is done exactly as in the previous use case.

An Additional Tool

If you have to manage a complex replication use case, you may want to check the open source project '[repmgr](https://repmgr.org/)'. It supports you to monitor the cluster of nodes or perform a failover.

Partitioning

If you have a table with a very huge amount of data, it may be helpful to scatter the data to different physical tables which share a common data structure. In such use cases, where DML statements concern only one of those physical tables, you can get great performance benefits from partitioning. Typically this is the case, if there is any timeline or a geographical distribution of the values of a column.

Declarative-partitioning-syntax: since version 10

Postgres 10 introduced a declarative partition-defining-syntax in addition to the previous table-inheritance-syntax. With this syntax the necessity to define an additional trigger disappears, but in comparison to the previous solution the functionality stays unchanged.

First, you define a master table containing a partitioning method which is `PARTITION BY RANGE (column_name)` in this example:

```
CREATE TABLE log (  
  id      int not null,  
  logdate date not null,  
  message varchar(500)  
) PARTITION BY RANGE (logdate);
```

Next, you create partitions with the same structure as the master and ensure, that only rows within the expected data range can be stored there. Those partitions are conventional, physical tables.

```
CREATE TABLE log_2015_01 PARTITION OF log FOR VALUES FROM ('2015-01-01') TO ('2015-02-01');  
CREATE TABLE log_2015_02 PARTITION OF log FOR VALUES FROM ('2015-02-01') TO ('2015-03-01');  
...  
CREATE TABLE log_2015_12 PARTITION OF log FOR VALUES FROM ('2015-12-01') TO ('2016-01-01');  
CREATE TABLE log_2016_01 PARTITION OF log FOR VALUES FROM ('2016-01-01') TO ('2016-02-01');  
...
```

Table-inheritance-syntax

First, you define a master table, which is a conventional table.

```
CREATE TABLE log (  
  id      int not null,  
  logdate date not null,  
  message varchar(500)  
);
```

Next, you create partitions with the same structure as the master table by using the table-inheritance mechanism `INHERITS (table_name)`. Additionally you must ensure that only rows within the expected data range can be stored in the derived tables.

```
CREATE TABLE log_2015_01 (CHECK (logdate >= DATE '2015-01-01' AND logdate < DATE '2015-02-01')) INHERITS (log);  
CREATE TABLE log_2015_02 (CHECK (logdate >= DATE '2015-02-01' AND logdate < DATE '2015-03-01')) INHERITS (log);  
...  
CREATE TABLE log_2015_12 (CHECK (logdate >= DATE '2015-12-01' AND logdate < DATE '2016-01-01')) INHERITS (log);  
CREATE TABLE log_2016_01 (CHECK (logdate >= DATE '2016-01-01' AND logdate < DATE '2016-02-01')) INHERITS (log);  
...
```

You need a function, which transfers rows into the appropriate partition.

```
CREATE OR REPLACE FUNCTION log_ins_function() RETURNS TRIGGER AS $$  
BEGIN  
  IF (NEW.logdate >= DATE '2015-01-01' AND NEW.logdate < DATE '2015-02-01' ) THEN  
    INSERT INTO log_2015_01 VALUES (NEW.*);  
  ELSIF (NEW.logdate >= DATE '2015-02-01' AND NEW.logdate < DATE '2015-03-01' ) THEN  
    INSERT INTO log_2015_02 VALUES (NEW.*);  
  ELSIF ...  
  ...  
  END IF;  
  RETURN NULL;  
END;  
$$  
LANGUAGE plpgsql;
```

The function is called by a trigger.

```
CREATE TRIGGER log_ins_trigger  
BEFORE INSERT ON log  
FOR EACH ROW EXECUTE PROCEDURE log_ins_function();
```

Further Steps

It's a good idea to create an index.

```
CREATE INDEX log_2015_01_idx ON log_2015_01 (logdate);
CREATE INDEX log_2015_02_idx ON log_2015_02 (logdate);
...
```

Many DML statements like `SELECT * FROM log WHERE logdate = '2015-01-15'`; act only on one partition and can ignore all the others. This is very helpful especially in such cases where a full table scan becomes necessary. The query optimizer has the chance to generate execution plans which avoid scanning unnecessary partitions.

In the shown example new rows will mainly go to the newest partition. After some years you can drop old partitions as a whole. This shall be done with the command `DROP TABLE` - not with a `DELETE` command. The `DROP TABLE` command is much faster than the `DELETE` command as it removes the complete partition in one single step instead of touching every single row.

Tablespace

The default behaviour of PostgreSQL is, that all data, indices, and management information is stored in subdirectories of a single directory. But this approach is not always suitable. In some situation you may want to change the storage area of one or more tables: data grows and may blow up partition limits, you may want to use fast devices like a ssd for heavily used tables, etc. . Therefore you need a technique to become more flexible.

Tablespaces offers the possibility to push data on arbitrary directories within your file system.

```
CREATE TABLESPACE fast LOCATION '/ssd1/postgresql/fastTablespace';
```

After the tablespace is defined it can be used in DDL statements.

```
CREATE TABLE t1(col_1 int) TABLESPACE fast;
```

Upgrade

When upgrading the PostgreSQL software, you must take care of the data in the cluster - depending on the question whether it is an upgrade of a major or a minor version. The PostgreSQL version number consists of two or three groups of digits, divided by colons. The first two groups denotes the major version and the third group (if present) denotes the minor version.

Upgrades within minor versions are simple. The internal data format does not change, so you only need to install the new software while the instance is down.

Upgrades of major versions may lead to incompatibilities of internal data structures. Therefore special actions may become necessary. There are several strategies to overcome the situation. In many cases upgrades of major versions additionally introduce some user-visible incompatibilities, so application programming changes might be required. You should read the release notes carefully.

pg_upgrade

`pg_upgrade` is a utility which modifies data files and system catalogs according to the needs of the new version. It has two major behaviors: In `--link` mode files are modified in place, otherwise the files are copied to a new location.

pg_dumpall

`pg_dumpall` is a standard utility to generate **logical** backups of the cluster. Files generated by `pg_dumpall` are plain text files and thus independent from all internal structures. When modifications of the data's internal structure become necessary (upgrade, different hardware architecture, different operating system, ...), such logical backups can be used for the data transfer from the old to the new system.

Replication

The Slony replication system offers the possibility to transfer data over different major versions. Using this, you can switch a replication slave to the new master within a very short time frame.

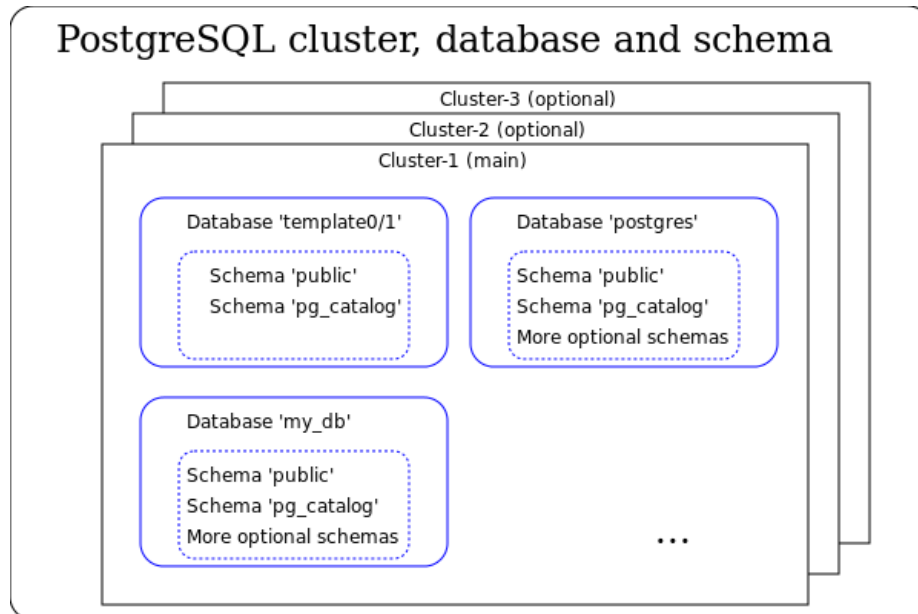
PostgreSQL offers replication in *logical streaming* format. With the actual version 9.5 this feature is restricted to the same versions of master and standby server, but it is planned to extend it for use in a heterogenous server landscape.

Terms

To promote a consistent use and understanding of important terms we list and define them here. For some terms we include short annotations to give a first introduction to the subject.

Database Cluster

Overview



Server (or Node)

A server is some (real or virtual) hardware where PostgreSQL is installed. In this document, the word *instance* is a different concept from *server*. See the definition of *instance* later in this document.

Cluster of Nodes

A set of nodes, which interchange information via replication.

Installation

After you have downloaded and installed PostgreSQL, you have a set of programs, scripts, configuration- and other files on a *server*. This set is called the 'Installation'. It includes all *instance* programs as well as some client programs like `psql`.

Server Database

The term *server database* is often used in the context of client/server connections to refer to an *instance* or a single *database*.

Cluster (or 'Database Cluster')

A cluster is a storage area (directory, subdirectories and files) in the file system, where a collection of databases plus meta-information resides. Within the database cluster there are also the definitions of global objects like users and their rights. They are known across the entire database cluster. (Access rights for an user may be limited to individual objects like a certain table or a certain schema. In that case, the user will not have this access rights to the other objects of the cluster.)

Within a database cluster there are at least three databases: 'template0', 'template1', 'postgres' and possibly more.

- 'template0': A template database, which may be used by the command `CREATE DATABASE` (template0 should never be modified)
- 'template1': A template database, which may be used by the command `CREATE DATABASE` (template1 may be modified by DBA)
- 'postgres': An empty database, mainly for maintenance purposes

Most PostgreSQL installations use only one database cluster. Its name is 'main'. But you can create more clusters on the same PostgreSQL installation, see tools `initdb` further down.

Instance (or 'Database Server Instance' or 'Database Server' or 'Backend')

An instance is a group of processes (on a UNIX server) or one service (on a Windows server) plus shared memory, which controls and manages exactly one *cluster*. Using IP terminology one can say that one instance occupies one IP/port combination, eg. the combination <http://localhost:5432>. It is possible that on a different port of the same *server* another instance is running. The processes (in a UNIX server), which build an instance, are called: postmaster (creates one 'postgres'-process per client-connection), logger, checkpointer, background writer, WAL writer, autovacuum launcher, archiver, stats collector. The role of each process is explained in the chapter [architecture](#).

If you have many *clusters* on your *server*, you can run many instances at the same machine - one per *cluster*.

Hint: Other publications sometimes use the term *server* to refer to an instance. As the term *server* is widely used to refer to real or virtual hardware, we do not use *server* as a synonym for *instance*.

Database

A database is a storage area in the file system, where a collection of objects is stored in files. The objects consist of data, metadata (table definitions, data types, constraints, views, ...) and other data like indices. Those objects are stored in the default database 'postgres' or in a newly created database.

The storage area for one database is organized as one subdirectory tree within the storage area of the database cluster. Thus a database cluster may contain multiple databases.

In a newly created database cluster (see below: `initdb`) there is an empty database with the name 'postgres'. In most cases this database stays empty and application data is stored in separate databases like 'finance' or 'engineering'. Nevertheless 'postgres' should not be dropped because some tools try to store temporary data within this database.

Schema

A schema is a namespace within a database: it contains named objects (tables, data types, functions, and operators) whose names can duplicate those of other objects existing in other schemas of this database. Every database contains the default schema 'public' and may contain more schemas. All objects of one schema must reside within the same database. Objects of different schemas within the same database may have the same name.

There is another special schema in each database. The schema 'pg_catalog' contains all system tables, built-in data types, functions, and operators. See also 'Search Path' below.

Search Path (or 'Schema Search Path')

A Search Path is a list of schema names. If applications use unqualified object names (e.g.: 'employee_table' for a table name), the search path is used to locate this object in the given sequence of schemas. The schema 'pg_catalog' is always the first part of the search path although it is not explicitly listed in the search path. This behaviour ensures that PostgreSQL finds the system objects.

initdb (OS command)

Despite of its name the utility `initdb` creates a new *cluster*, which contains the 3 *databases* 'template0', 'template1' and 'postgres'.

createdb (OS command)

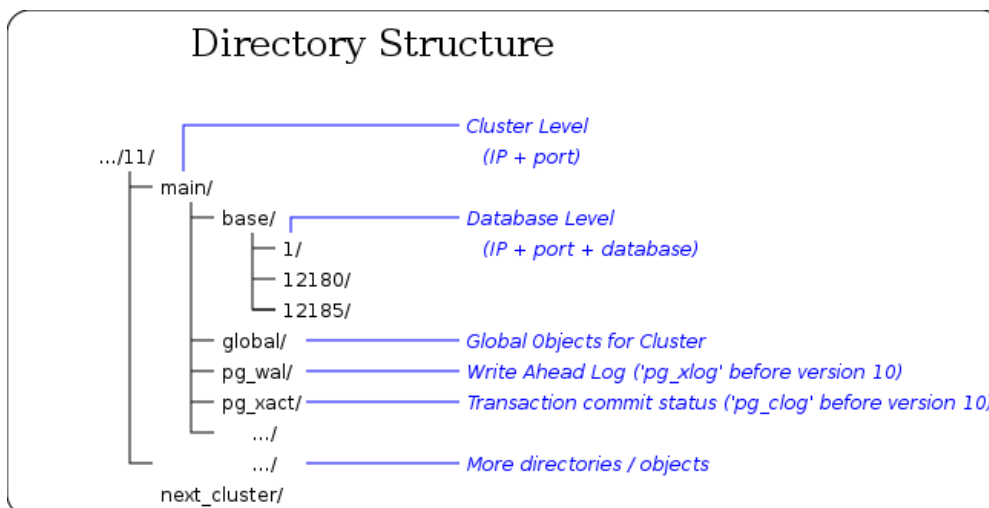
The utility `createdb` creates a new *database* within the actual *cluster*.

CREATE DATABASE (SQL command)

The SQL command `CREATE DATABASE` creates a new *database* within the actual *cluster*.

Directory Structure

A *cluster* and its *databases* consists of files, which hold data, actual status information, modification information and a lot more. Those files are organized in a fixed way under one directory node.



Consistent Writes

Shared Buffers

Shared buffers are RAM pages, which mirror pages of data files on disc. They exist due to performance reasons. The term *shared* results from the fact that a lot of processes read and write to that area.

'Dirty' Page

Pages in the *shared buffers* mirror pages of data files on disc. When clients request changes of data, those pages get changed without - provisionally - a change of the related pages on disc. Until the *background writer* writes those modified pages to disc, they are called 'dirty' pages.

Checkpoint

A checkpoint is a special point in time where it is guaranteed that the database files are in a consistent state. At checkpoint time all change records are flushed to the WAL file, all dirty data pages (in shared buffers) are flushed to disc, and at last a special checkpoint record is written to the WAL file.

The instance's checkpointer process automatically triggers checkpoints on a regular basis. Additionally they can be forced by issuing the command CHECKPOINT in a client program. For the database system it takes a lot of time to perform a checkpoint - because of the physical writes to disc.

WAL File

WAL files contain the changes which are applied to the data by modifying commands like INSERT, UPDATE, DELETE or CREATE TABLE This is redundant information as it is also recorded in the data files (for better performance at a later time). According to the configuration of the instance there may be more information within WAL files. WAL files reside in the pg_wal directory (which was named pg_xlog before version 10), have a binary format and a fixed size of 16MB. When they are no longer needed, they get recycled by renaming and reusing their already allocated space.

A single information unit within a WAL file is called a *log record*.

Hint: In the PostgreSQL documentation and in related documents there are a lot of other, similar terms which refer to what we denote as *WAL file* in this Wikibook: segment, WAL segment, logfile (don't mix it with the term *logfile*, see below), WAL log file,

Logfile

The *instance* logs and reports warning and error messages about special situations in readable text files. These logfiles can reside at any place in the directory structure of the *server* and are not part of the *cluster*.

Hint: The term 'logfile' does not relate to the other terms of this subchapter. It is mentioned here because the term sometimes is used as a synonym for what we call *WAL file* - see above.

Log Record

A log record is a single information unit within a *WAL file*.

Segment

The term *segment* is sometimes used as a synonym for *WAL file*.

MVCC

Multiversion Concurrency Control (MVCC) is a common database technique to accomplish two goals: First, it allows the management of parallel running transactions on a logical level and second, it ensures high performance for concurrent read and write actions. It is implemented as follows: Whenever some values of an existing row change, PostgreSQL writes a new version of this row to the database without deleting the old one. In such situations the database contains multiple versions of the row. In addition to their regular data the rows contain transaction IDs which allows to decide, which other transactions will see the new or the old row. Hence other transactions sees only those values (of other transactions), which are committed.

Outdated old rows are deleted at a later time by the utility `vacuumdb` respectively the SQL command `vacuum`.

Backup and Recovery

The term *cold* as an addition to the backup method name indicates that with this method the instance must be stopped to create a useful backup. In contrast, the addition 'hot' denotes methods where the instance **MUST** run (and hence changes to the data may occur during backup actions).

(Cold) Backup (file system tools) (<http://www.postgresql.org/docs/current/static/creating-cluster.html>)

A cold backup is a consistent copy of all files of the *cluster* with OS tools like `cp` or `tar`. During the creation of a cold backup the *instance* must **not** run - otherwise the backup is useless. Hence you need a period of time in which applications do not use any *database* of the *cluster* - a continuous 7×24 operation mode is not possible. And secondly: the cold backup works only on the *cluster* level, not on any finer granularity like *database* or *table*.

Hint: A cold backup is sometimes called an "offline backup".

(Hot) Logical Backup (pg_dump utility) (<http://www.postgresql.org/docs/current/static/backup-dump.html>)

A logical backup is a consistent copy of the data within a *database* or some of its parts. It is created with the utility `pg_dump`. Although `pg_dump` may run in parallel to applications (the *instance* must be up), it creates a consistent snapshot as of the time of its start. `pg_dump` supports two output formats. The first one is a text format containing SQL commands like `CREATE` and `INSERT`. Files created in this format may be used by `psql` to restore the backed-up data. The second format is a binary format and is called the 'archive format'. Files with this format can be used to restore its data with the tool `pg_restore`.

As mentioned, `pg_dump` works at the *database* level or smaller parts of *databases* like *tables*. If you want to refer to the *cluster* level, you must use `pg_dumpall`. Please notice, that important objects like *users/roles* and their rights are always defined at *cluster* level.

Hint: A logical backup is one form of an "online backup".

(Hot) Physical Backup or 'Base Backup' (<http://www.postgresql.org/docs/current/static/continuous-archiving.html>)

A physical backup is a **possibly inconsistent** copy of the files of a *cluster*, created with an operating system utility like `cp` or `tar`. At first glance such a backup seems to be useless. To understand its purpose, you must know PostgreSQL's recover-from-crash strategy.

At all times and independent from any backup/recovery action, PostgreSQL maintains *WAL files* - primarily for crash-safety purposes. *WAL files* contain *log records*, which reflect all changes made to the data. In the case of a system crash those *log records* are used to recover the *cluster* to a consistent state. The recover process searches the timestamp of the last *checkpoint* and replays all subsequent *log records* in chronological order against this version of the *cluster*.

Through this action the *cluster* returns to a consistent state and will contain all changes up to the last `COMMIT`.

The existence of a physical backup plus *WAL files* in combination with this recovery-from-crash technique can be used for backup/recovery purposes. To implement this, you need a physical backup (which may reflect an inconsistent state of the *cluster*) and which acts as the starting point. Additionally you need all *WAL files* since the point in time when you have created this backup. The recover process uses the described recovery-from-crash technique and replays all *log records* in the *WAL files* against the backup. In the exact same way as before, the *cluster*

comes to a consistent state and contains all changes up to the last COMMIT.

Please keep in mind, that physical backups work only on cluster level, not on any finer granularity like database or table.

Hint: A physical backup is one form of an "online backup".

PITR: Point in Time Recovery (<http://www.postgresql.org/docs/current/static/continuous-archiving.html>)

If the previously mentioned *physical backup* is used during recovery, the recovery process is not forced to run up to the latest available timestamp. Via a parameter you can stop it at a time in the past. This leads to a state of the *cluster* at this moment. Using this technique you can restore your *cluster* to a time, which is between the time of creating the *physical backup* and the end of the last *WAL file*.

Standalone (Hot) Backup

The standalone backup is a special variant of the physical backup. It offers online backup (the *instance* keeps running) but it lacks the possibility of *PITR*. The recovery process recovers always up to the end of the standalone backup process. *WAL files*, which arise after this point in time, cannot be applied to this kind of backup. Details are described here (<http://www.postgresql.org/docs/current/static/continuous-archiving.html>).

Archiving

Archiving is the process of copying *WAL files* to a failsafe location. When you plan to use *PITR* you must ensure that the sequence of *WAL files* is saved for a longer period. To support the process of copying *WAL files* at the right moment (when they are completely filled and a switch to the next *WAL file* has taken place), PostgreSQL runs the *archiving process* which is part of the *instance*. This process copies *WAL files* to a configurable destination.

Recovering

Recovering is the process of playing *WAL files* against a *physical backup*. One of the involved steps is the copy of the *WAL files* from the failsafe archive location to its original location in '*pg_xlog*'. The aim of recovery is bringing the *cluster* into a consistent state at a defined timestamp.

Archive Recovery Mode

When recovering takes place, the *instance* is in *archive recovery mode*.

Restartpoint

A restartpoint is an action similar to a *checkpoint*. Restartpoints are only performed when the instance is in *archive recovery mode* or in *standby mode*.

Timeline

After a successful recovery PostgreSQL transfers the *cluster* into a new timeline to avoid problems, which may occur when *PITR* is reset and *WAL files* reapplied (e.g.: to a different timestamp). Timeline names are sequential numbers: 1, 2, 3,

Replication

Replication is a technique to send data, which was written within a *master server*, to one or more *standby servers* or even another *master server*.

Master Server

The master server is an *instance* on a *server* which sends data to other *instances* in addition to its local processing of data.

Standby Server

The standby server is an *instance* on a *server* which receives information from a *master server* about changes of its data.

Warm Standby Server

A warm standby server is a running *instance*, which is in *standby_mode* (recovery.conf file). It continuously reads and processes incoming *WAL files* (in the case of *log-shipping*) or *log records* (in the case of *streaming replication*). It does not accept client connections.

Hot Standby Server

A hot standby server is a warm standby server with the additional flag *hot_standby* in postgres.conf. It accepts client connections and read-only queries.

Synchronous Replication

Replication is called *synchronous*, when the *standby server* processes the received data immediately, sends a confirmation record to the *master server* and the *master server* delays its COMMIT action until he has received the confirmation of the *standby server*.

Asynchronous Replication

Replication is called *asynchronous*, when the *master server* sends data to the *standby server* and does not expect any feedback about this action.

Streaming Replication

The term is used when *log entries* are transferred from *master server* to *standby server* over a TCP connection - in addition to their transfer to the local *WAL file*. Streaming replication is *asynchronous* by default but can also be *synchronous*.

Log-Shipping Replication

Log shipping is the process of transferring *WAL files* from a *master server* to a *standby server*. Log shipping is an asynchronous operation.

Extensions

PostgreSQL offers an extensibility architecture and implements its internal datatypes, operators, functions, indexes, and more on top of it. This architecture is open for everybody to implement and add his own functionality to the PostgreSQL system. You can define new datatypes with or without special operators and functions as needed by your use case. After you have added them, you have the best of two worlds: the special functionalities you have created plus the standard functionality of a database system like ACID, SQL, security, standard data types, WAL, client APIs, An introduction to extensibility is given in the [PostgreSQL documentation \(https://www.postgresql.org/docs/current/extend.html\)](https://www.postgresql.org/docs/current/extend.html).

Over time the community has developed a set of extensions which are useful for their own needs and for a great number of applications - sometimes even for the requirements and definitions given by standardization organisations. Some popular examples are

- Data types, operators, and function for the handling of **spatial data** like points, polylines, overlaps(), ... as defined by [OSGeo](#) and [SQL Multimedia and Application Packages Part 3: Spatial](#).
- Functionality for **full text** search as defined by [SQL Multimedia and Application Packages Part 2: Full-Text](#).
- Access to **data outside** the current database (other PostgreSQL instance, other SQL, NoSQL or BigData database system, LDAP, flat files like csv, json, xml) as defined by [SQL Part 9: Management of External Data](#).

The lifecycle of such an extension starts with the implementation of its features by a group of persons or a company. After publishing, the extension may be used and further expanded by other persons or companies of the community. Sometimes such extensions keep independent from the PostgreSQL system, e.g.: PostGIS, in other cases they are delivered with the standard download and explicitly listed as an [Additional Supplied Module \(https://www.postgresql.org/docs/current/contrib.html\)](#) within the documentation with hints how to install them. And in rare cases extensions are incorporated into the core system so that they become a native part of PostgreSQL.

To activate and use an extension, you must download and install the necessary files (if not delivered with the standard download) and issue the command `CREATE EXTENSION <extension_name>;` within an SQL client like `psql`. To control which extensions are already installed use: `\dx` within `psql`.

PostGIS

PostGIS is a project which extends PostgreSQL with a rich set of 2D and 3D spacial data types plus associated functions, operators and index types as defined by [OSGeo](#) and [SQL Multimedia and Application Packages Part 3: Spatial](#). Typically data types are *polygon* or *multipoint*, typical functions are `st_length()` or `st_contains()`. The appropriated index type for spatial objects is the [GiST index](#).

The PostGIS project has its own [representation on the WEB \(http://postgis.net/\)](http://postgis.net/) where all its aspects are described, especially the download process and the activation of the extension itself.

Foreign Data Wrappers

Foreign Data Wrappers (FDW) are PostgreSQL extensions which offers access to data outside of the actual database and instance. There are different types of data wrappers:

- One wrapper to other PostgreSQL instances: `postgres_fdw`
- A lot of wrappers to other relational database systems like Oracle, MySQL, MS SQL Server, ...
- A lot of wrappers to NoSQL database systems: CouchDB, MongoDB, Cassandra, ...
- Generic wrappers to ODBC and JDBC
- A lot of wrapper to files of different formats: csv, xml, json, tar, zip, ... (`file_fdw`)
- LDAP wrapper
- ... and more.

A comprehensive [list \(https://wiki.postgresql.org/wiki/Foreign_data_wrappers\)](https://wiki.postgresql.org/wiki/Foreign_data_wrappers) gives an overview.

The technique of FDW is defined in the SQL standard [Part 9: Management of External Data](#).

Here is an example how to access another PostgreSQL instance via FDW.

```
-- Install the extension to other PostgreSQL instances
CREATE EXTENSION postgres_fdw;

-- Define the connection to a database/instance at a different server
CREATE SERVER remote_geo_server
  FOREIGN DATA WRAPPER postgres_fdw
  OPTIONS (host '10.10.10.10', port '5432', dbname 'geo_data');

-- Define a user for the connection (The remote user must have access rights at the remote database)
CREATE USER MAPPING FOR CURRENT_USER
  SERVER remote_geo_server
  OPTIONS (user 'geo_data_user', password 'xxx');

-- Define two foreign tables via an IMPORT command ...
```

```
IMPORT FOREIGN SCHEMA geo_schema
LIMIT TO (city, point_of_interest)
FROM SERVER remote_geo_server
INTO my_schema;

-- .. and another foreign table via an explicit definition
CREATE FOREIGN TABLE remote_person (
  id          SERIAL,
  person_name TEXT          NOT NULL,
  city_id     INT4          NOT NULL
)
SERVER remote_geo_server
OPTIONS(schema_name 'geo_schema', table_name 'person');
```

After the execution of the above statements you have access to the three tables `city`, `point_of_interest` and `remote_person` with the usual DML commands `SELECT`, `UPDATE`, `COMMIT`, Nevertheless the data keeps at the 'remote' server (10.10.10.10), queries are executed there, and only the results of queries are transferred via network to the actual instance and your client application.

```
SELECT count(*) FROM city; -- table 'city' resides on a different server
```

Bidirectional Replication (BDR)

BDR is an extension which allows replication in both directions between involved (master-) nodes in parallel to their regular read and write activities of their client applications. So it realizes a multi-master replication. Actually the project is a [standalone project \(http://bdr-project.org/docs/next/index.html\)](http://bdr-project.org/docs/next/index.html). But multiple technologies emerging from BDR development have already become an integral part of core PostgreSQL, such as [Event Triggers \(https://www.postgresql.org/docs/current/event-triggers.html\)](https://www.postgresql.org/docs/current/event-triggers.html), [Logical Decoding \(https://www.postgresql.org/docs/current/logicaldecoding.html\)](https://www.postgresql.org/docs/current/logicaldecoding.html), [Replication Slots \(https://www.postgresql.org/docs/current/logicaldecoding-explanation.html#LOGICALDECODING-REPLICATION-SLOTS\)](https://www.postgresql.org/docs/current/logicaldecoding-explanation.html#LOGICALDECODING-REPLICATION-SLOTS), [Background Workers \(https://www.postgresql.org/docs/current/bgworker.html\)](https://www.postgresql.org/docs/current/bgworker.html), and more.

Retrieved from "https://en.wikibooks.org/w/index.php?title=PostgreSQL/Print_version&oldid=3836743"

This page was last edited on 18 May 2021, at 03:13.

Text is available under the Creative Commons Attribution-ShareAlike License.; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy.