

# Introduction to PostgreSQL

Ryan Booz & Grant Fritchey



# Ryan Booz

## Redgate Advocate

 [@ryanbooz](https://twitter.com/@ryanbooz)

 [/in/ryanbooz](https://www.linkedin.com/in/ryanbooz/)

 [www.softwareandbooz.com](http://www.softwareandbooz.com)

 [youtube.com/@ryanbooz](https://youtube.com/@ryanbooz)



# Grant Fritchey

## DevOps Advocate

Microsoft PostgreSQL MVP  
AWS Community Builder

 @gfrtchey

 Grant.Fritchey@Red-Gate.com

 /in/gfrtchey

 scarydba.com



# Why PostgreSQL?

# Licensing

- Extremely permissive, open-source license
- Core features can never be paywalled
- Most extensions piggyback on the PostgreSQL License





# Amazon RDS

This Photo by Unknown Author is licensed under [CC BY](#)



# Extensibility



- >30 hooks available to extensions
- Any supported language, although SQL, C, and Rust are most popular

# Extensibility



- Not just for extending functionality
- Package commonly used SQL functions/procedures
- Versioning required and helpful

# Lay the Foundation

# **Instance = Cluster**

SQL Server

PostgreSQL

# User = Role

SQL Server

PostgreSQL

# Transaction Log ≈ Write Ahead Log

SQL Server

PostgreSQL

# **Row = Row ∈ Tuple**

SQL Server

PostgreSQL

PostgreSQL

**[Table One] = "Table One"**

SQL Server

PostgreSQL

# BULK INSERT ≈ COPY

SQL Server

PostgreSQL

# TOAST

The Oversized-Attribute Storage Technique

# Transaction Isolation

- Read Committed
- Repeatable Read
- Serializable Read
- **Read Uncommitted (not supported)**

# PostgreSQL MVCC at a glance

- Readers don't block writers
- Tuples maintain on-row transaction visibility (xmin, xmax)
- UPDATES always create new rows
  - Heap-only Tuples (HOT) attempts to mitigate some growth
- Popular MVCC talk: [https://www.youtube.com/watch?v=gAE\\_MSQtqnQ](https://www.youtube.com/watch?v=gAE_MSQtqnQ)

# MVCC at a glance

- Readers don't block writers
- Tuples maintain on-row transaction visibility values (xmin, xmax)
- **UPDATES always create new rows**
  - Heap-only Tuples (HOT) attempts to mitigate some growth
  - Popular MVCC talk: [https://www.youtube.com/watch?v=gAE\\_MSQtqnQ](https://www.youtube.com/watch?v=gAE_MSQtqnQ)

# Case Rules

- Case sensitive string comparison (LIKE/ILIKE)
- All object names are coerced to lower case
- "Double quote" to qualify names
- Prefer lower\_snake\_case

# SQL Things

# Case Rules

SQL Server

```
SELECT [Name]  
      FROM [Purchasing].[Vendor]  
     WHERE BusinessEntityID = 1492
```

Postgres

```
SELECT "Name"  
      FROM "Purchasing"."Vendor"  
     WHERE "BusinessEntityID" = 1492;
```

# Case Rules

SQL Server

```
SELECT [Name]  
      FROM [Purchasing].[Vendor]  
     WHERE BusinessEntityID = 1492
```

Postgres

```
SELECT name  
      FROM purchasing.vendor  
     WHERE business_entity_id = 1492;
```

# LIMIT ≈ TOP

SQL Server

```
SELECT TOP(10) *
    FROM [Purchasing].[Vendor]
```

Postgres

```
SELECT *
    FROM purchasing.vendor
    LIMIT 10;
```

# LIMIT/OFFSET ≈ OFFSET/FETCH

SQL Server

```
SELECT *
    FROM [Purchasing].[Vendor]
  ORDER BY name
OFFSET 10 ROWS
FETCH 10 ROWS ONLY
```

Postgres

```
SELECT *
    FROM purchasing.vendor
  ORDER BY name
LIMIT 10 OFFSET 10;
```

# ORDER/GROUP BY

SQL Server

```
SELECT SUBSTRING(name,1,1) nameCohort, SUM(amount) totalSales  
FROM [Purchasing].[Sales]  
GROUP BY SUBSTRING(name,1,1)  
ORDER BY nameCohort, SUM(amount)
```

Postgres

```
SELECT SUBSTRING(name,1,1) name_cohort, SUM(amount) total_sales  
FROM purchasing.sales  
GROUP BY name_cohort  
ORDER BY name_cohort, total_sales;
```

# ORDER/GROUP BY

SQL Server

```
SELECT SUBSTRING(name,1,1) nameCohort, SUM(amount) totalSales  
FROM [Purchasing].[Sales]  
GROUP BY SUBSTRING(name,1,1)  
ORDER BY nameCohort, SUM(amount)
```

Postgres

```
SELECT SUBSTRING(name,1,1) name_cohort, SUM(amount) total_sales  
FROM purchasing.sales  
GROUP BY 1  
ORDER BY 1, 2;
```

# DAY/MONTH/YEAR() ≈ DATE\_PART()

SQL Server

```
SELECT COUNT(*), YEAR(OrderDate)
FROM Sales.Invoices
GROUP BY YEAR(OrderDate)
ORDER BY YEAR(OrderDate) DESC
```

Postgres

```
SELECT COUNT(*), DATE_PART('year',order_date) order_year
FROM sales.invoices
GROUP BY order_year
ORDER BY order_year DESC;
```

# **SYSDATETIME() ≈ NOW()**

SQL Server

```
SELECT SYSDATETIME();
```

Postgres

```
SELECT now();
```

# Date Math & Intervals

SQL Server

```
SELECT COUNT(*), YEAR(OrderDate)
FROM Sales.Invoices
WHERE OrderDate > DATEADD(MONTH,-1,SYSDATETIME())
GROUP BY YEAR(OrderDate)
```

Postgres

```
SELECT COUNT(*), DATE_PART('year',order_date) order_year
FROM sales.invoices
WHERE order_date > NOW() - INTERVAL '1 MONTH'
GROUP BY order_year;
```

# LATERAL ≈ APPLY

- CROSS/OUTER APPLY in T-SQL
  - Often used as an “optimization fence”
  - Outer query iterates an inner query
  - Helpful when JOIN isn’t possible
  - Retrieve top ordered row for each item in a list (ie. “most recent value for each item”)
  - Iterating a function that references the outer query
- LATERAL JOIN allows the similar functionality

# LATERAL ≈ APPLY

SQL Server

```
SELECT id, name, dayMonth, total
      FROM Vendor v1
CROSS APPLY (
    SELECT TOP(1) dayMonth, SUM(amount) total
      FROM Sales
     WHERE customerID = v1.id
     GROUP BY dayMonth
     ORDER BY SUM(amount) DESC
) s1
```

# LATERAL ≈ APPLY

Postgres

```
SELECT id, name, day_month, total
  FROM vendor v1
INNER JOIN LATERAL (
    SELECT day_month, SUM(amount) total
      FROM sales
     WHERE customer_id = v1.id
     GROUP BY day_month
     ORDER BY total DESC
     LIMIT 1
) on true s1;
```

# Anonymous Code blocks

- Default query language is SQL
- Does not support anonymous code blocks
  - Variables
  - Loops
  - T-SQL like features
- Requires DO blocks, Functions, or SPROCs
  - Default language is pl/pgsql ≈ T-SQL

# Inline T-SQL vs pl/pgsql

SQL Server

```
DECLARE @PurchaseName AS NVARCHAR(50)  
SELECT @PurchaseName = [Name]  
    FROM [Purchasing].[Vendor]  
    WHERE BusinessEntityID = 1492  
PRINT @PurchaseName
```

# Inline T-SQL vs pl/pgsql

Postgres

```
DO $$  
    DECLARE purchase_name TEXT;  
  
BEGIN  
    SELECT name  
    FROM purchasing.vendor  
    WHERE business_entity_id = 1492 INTO purchase_name;  
  
    RAISE NOTICE '%', purchase_name;  
END;  
$$
```

# Functions

## SQL Server

- Scaler
- Inline Table Valued
- Multi-statement Table Valued
- No tuning hints for query planner
- TSQL Only

## Postgres

- Scaler
- Table/Set Types
- Triggers
- Query planner tuning parameters
- Any procedural language, but typically PL/pgSQL

# Stored Procedures

## SQL Server

- Complex logic
- Multi-language
- Contained in calling transaction scope
- With some work, can return results

## Postgres

- Complex logic
- Multi-language
- Can issue COMMIT/ROLLBACK and keep processing
- Can return a single INOUT value
- Introduced in PostgreSQL 11

# Triggers

- Functions applied to tables
- INSERT/UPDATE/DELETE
- BEFORE/AFTER/INSTEAD OF
- ROW and STATEMENT
- Conditional
- NEW and OLD internal variables
- ...Reusable across tables!

# Create Trigger Function

Postgres

```
CREATE OR REPLACE FUNCTION user_log()
RETURNS TRIGGER AS
$BODY$
BEGIN
    INSERT INTO UserLog (Username, Message) VALUES (NEW.Username, NEW.Message)
END;
$BODY$
LANGUAGE plpgsql;
```

# Apply Trigger

Postgres

```
CREATE TRIGGER tu_users
AFTER UPDATE
ON Users
FOR EACH ROW
WHEN (OLD.FirstName IS DISTINCT FROM NEW.FirstName)
EXECUTE PROCEDURE user_log();
```

**DEMO!**



PASS Session

# The Accidental PostgreSQL Database Administrator



Wednesday, Nov. 6



5:00-6:00



Room 447



PASS Session

# MSSQL and PostgreSQL

## Similarities and Differences for the SQL Server DBA



Friday, Nov. 1



1:00-2:00



Room 445-446



# Installation

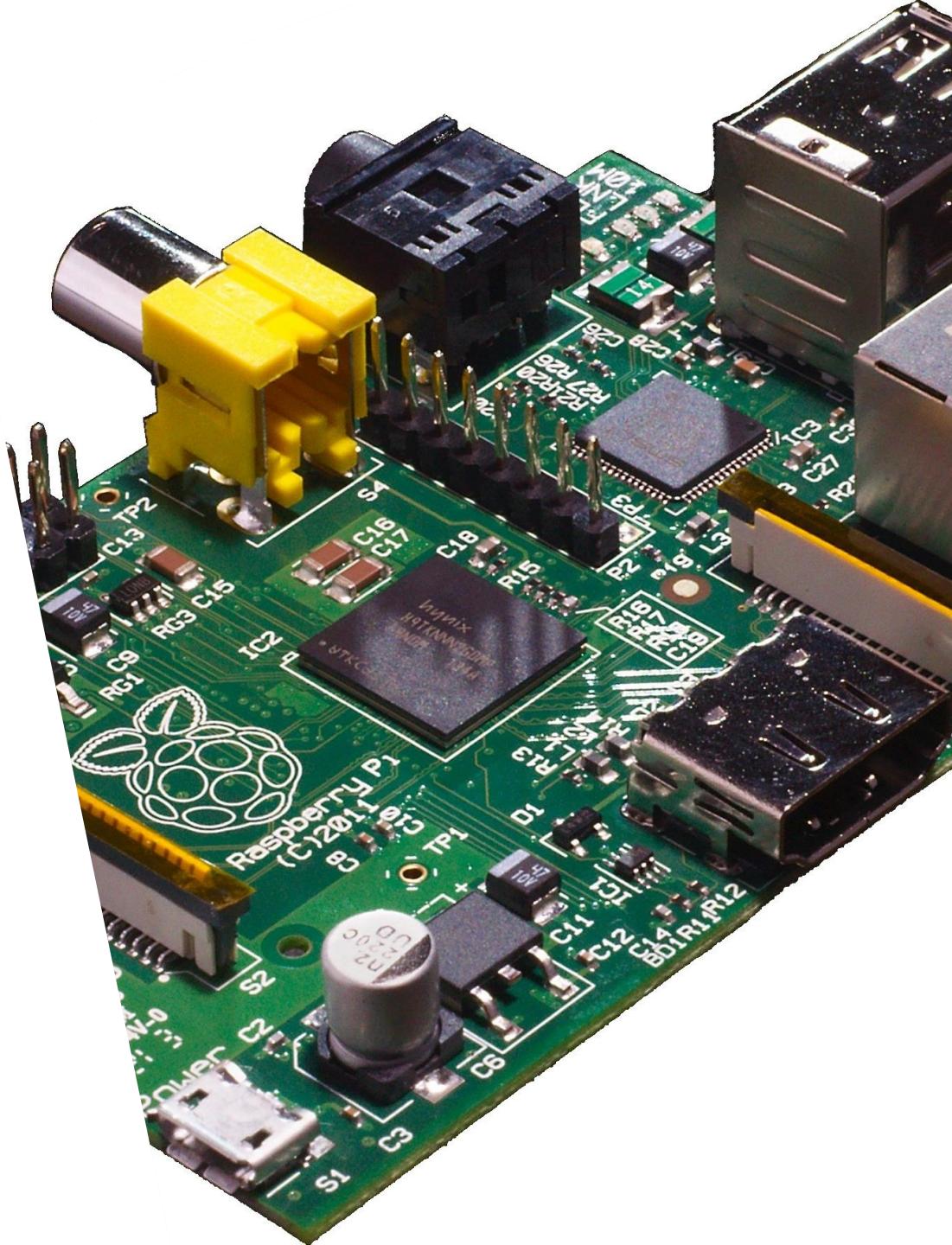
# Supported Operating Systems



- Linux
- Windows
- Containers
- Cloud

# Minimum Hardware

- Chip sets
  - x86 family
  - ARM
  - SPARC
  - PowerPC
  - S/390
  - MIPS
  - RISC-V
  - PA-RISC
- Recommended Minimums
  - > 1ghz speed CPU
  - > 2gb memory
  - > 50kb storage



# Packaged Bits

- Core
  - For the “hackers”
- Packages
  - OS
  - Special functions



# Containers

- Easiest way to get started
- Get in the habit of setting up volumes
- Multiple forks available



**DEMO!**



# PostgreSQL Tools

# psql

- Always available (almost)
- Command line
- Can work remote

```
shell 7.4.5
C:\Users\grant> psql -?
psql is the PostgreSQL interactive terminal.

Usage:
  psql [OPTION]... [DBNAME [USERNAME]]

General options:
  -c, --command=COMMAND      run only single command
  -d, --dbname=DBNAME        database name to connect to
  -f, --file=FILENAME         execute commands from file
  -l, --list                  list available databases
  -v, --set=, --variable=NAME=VALUE
                               set psql variable NAME to VALUE
                               (e.g., -v ON_ERROR_STOP=1)
  --version                  output version information
  --no-psqlrc                do not read startup file
  ("one"), --single-transaction
                               execute as a single transaction
  --help[=options]
  --help=commands
  --help=variables

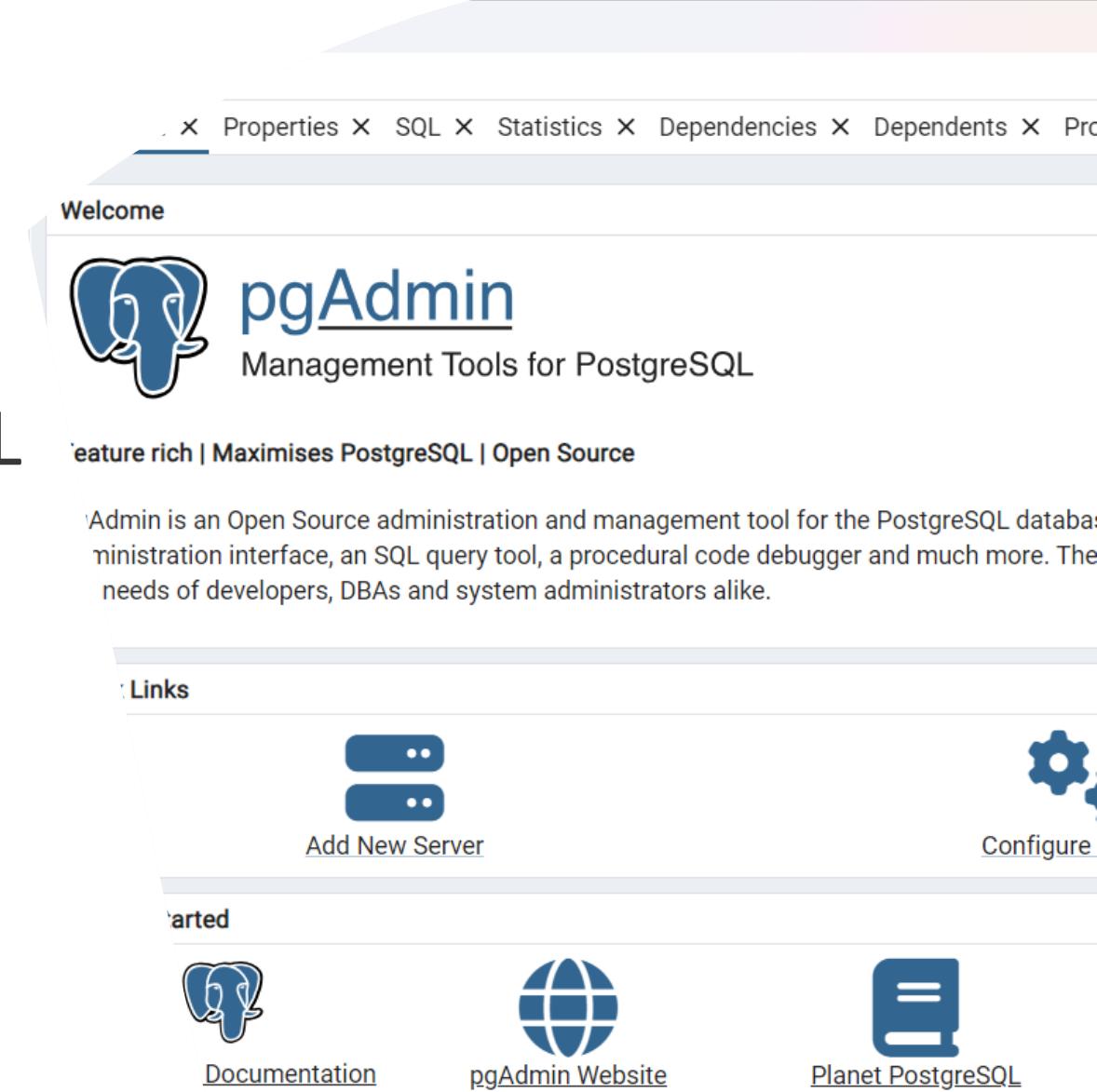
Output options:
  -echo-all
  -echo-errors
  -echo-queries
  -echo-hidden
  -log-file=FILENAME
  -readline
```

# pgAdmin

- Unofficial, "Official" PostgreSQL

GUI

- Not installed by default
- Web version for remote work
- Extensions



# Azure Data Studio/VS Code

- PostgreSQL Plugin
- Also works in VSCode
- Good for basics
- Good integration with source control



Azure Data Stu

New ▾

Open file...

Open fo

Create a connection  
Connect to a database instance  
through the connection dialog.

# DBeaver

- Open source
  - Also paid
- Multi-platform
  - Good and bad
- Extensions
- Some code completion



**DEMO!**



# Databases

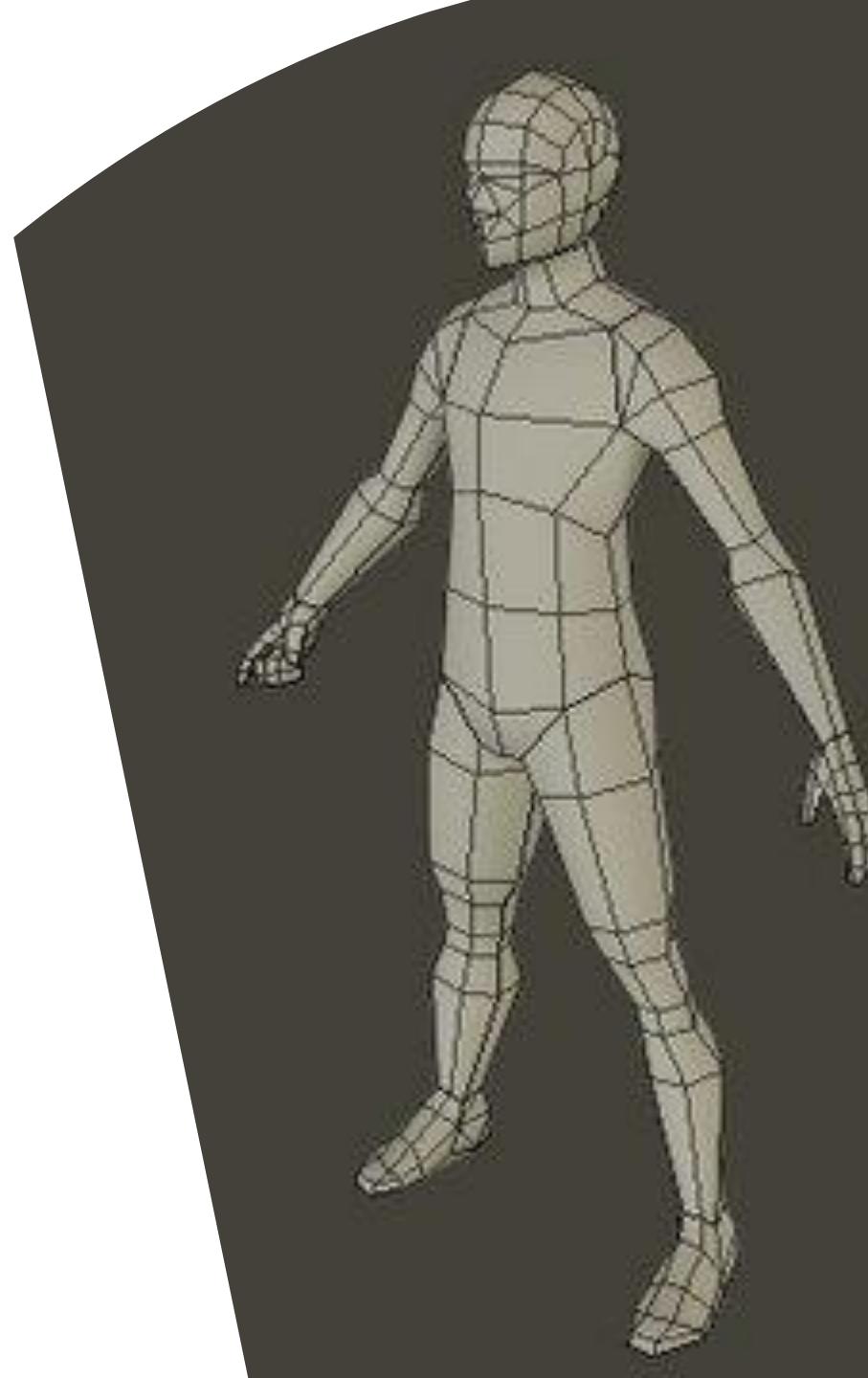
# What Are Databases?

- Primary means of organizing storage, access and behavior
- Equivalent to SQL Server databases
- Control storage
- 3 by default
  - postgres
  - template0
  - template1



# Templates

- Think, model db
- template0 – leave this alone
- template1 – customizable
- Custom templates
- You can choose the template



# CREATE DATABASE

- Defaults to template1
- You can pick the template
- You can define the owner



# ALTER DATABASE

- Change database behaviors
- Change database properties
- Modify behaviors of code and objects within database



# DROP DATABASE

- Boom!
- IF EXISTS is a good idea
- WITH (FORCE) is probably not



**DEMO!**



# Restores and Backups

**“You are only as good as your last restore”**

- Kimberly Tripp

# Backups



- pg\_dump
- pg\_dumpall
- pg\_basebackup
- Restore from dump
- pg\_restore

# pg\_dump



- An export, not a true backup
- Not transactionally aware
- No point in time recovery
- One database at a time
- pg\_dumpall

# Common pg\_dump Behaviors



- -F – provides a file format, common to compress in tar file
- -c – adds drop commands so you don't have to drop & recreate the db (sometimes)
- -a – data only
- -s – schema only
- -C – creates a database from dump

# pg\_basebackup



- Backs up the server
- Starting point for point in time recovery

# WAL Archiving



- Write Ahead Log (WAL)
- Log backups for point in time recovery

# WAL Archiving



- WAL level
  - Replica
  - Logical
- Archive mode enabled
- Define archive command
  - File copy, but more

**<https://github.com/ryanbooz/bluebox>**

- [GitHub - ryanbooz/bluebox: An updated sample database for PostgreSQL, building off of the Pagila database](#)

# Restore



- “Replay” a dump
- pg\_restore

# Restore to Point In Time



1. Run `pg_switch_wal`
2. Stop the server
3. Clean data directory or make a new one
4. Restore the Base Backup files
5. Remove all files from `pg_wal/`
6. Edit `postgresql.conf` for restore command and point in time
7. Create a file called `recovery.signal` in the cluster data directory
8. Start the service

**DEMO!**



# Database Objects

# Data Types

- Twenty (20) Categories of data types
- Depending on how you count,  
~125 data types
- We won't be listing them all



# Common Data Types

- Numeric
  - Integers, decimals & math
- Character
  - No need for Unicode
- Date/Time
- Binary
- Enums
- Geometric
- JSON
- Array



# Data Type Behaviors

- Fixed or variable length,  
depending on type
- Most can be NULL
- DEFAULT values can be defined



# Tables and Columns

- CREATE TABLE
- Columns have defined data types
- SERIAL vs. IDENTITY
- Belongs to a schema
- Owned by the creator, unless specified otherwise



# Constraints

- NOT NULL
- UNIQUE
- Primary key
- Foreign key
- Check
- Exclusion



# Primary Key

- Almost identical to a unique constraint
- No NULL values
- Defined by column, or by altering table



# Foreign Key

- Created against primary key or unique constraint
- Enforces referential integrity
  - But NULL values can be equal
- ON DELETE RESTRICT
- CASCADE



# Check Constraints

- Enforced on FALSE evaluation
  - Which means NULL comparisons count as TRUE
- Created with columns or ALTER TABLE



# Exclusion Constraints

- Similar to UNIQUE
- Involves more complex comparisons
- A strict equality comparison would be the same as UNIQUE



# Indexes

- PostgreSQL tables are heaps
- No concept of a clustered index
- CREATE and DROP  
CONCURRENTLY
- Must be table owner to add  
indexes



# Index Types

- B-Tree
- Hash
- GiST
- SP-GiST
- GIN
- BRIN



# Views

- Just a query
- Treated as a table in SQL
- Supports materialized views



**DEMO!**



# Procedural Languages

# Trusted vs. Untrusted Languages



- Trusted languages preserve security by preventing certain operations
  - DB internal processes
  - OS-level access
  - Non-privileged users can safely use

# Trusted vs. Untrusted Languages



- Untrusted languages do not provide security safeguards
- Creator is responsible for security
- Only superusers can create functions in untrusted languages

# Core Procedural PostgreSQL Languages



- SQL\* (T)
- PL/pgSQL (T)
- PL/Python (U)
- PL/Tcl (T/U)
- PL/Perl (T/U)

# Languages Available by Extension



- PL/Java (T/U)
- pl/dotnet (U)
- PL/Rust (T/U)
- PL/R (U)
- PI/v8js (T)

# Core Functionality of Procedural Languages



- Blocks
- Variables
- Error handling
- If/Then
- Case
- Functions
- Procedures

# Functions



- Implemented for returning data from PostgreSQL database
- Similar to Oracle functions or SQL Server procedures

# Function Basics



- Use CREATE FUNCTION command
- Provide a name
- Define result set, value or record/tuple
- Create a code block
- Define language used

# Function Behaviors



- Code blocks
  - Dollar quoting, \$\$
- Parameters
- Variables

# Function Example

```
CREATE FUNCTION genre_list ()
RETURNS TABLE (genre_id int, genre_name TEXT)
AS $$

BEGIN
    RETURN QUERY
    SELECT
        fg.genre_id,
        fg.name
    FROM
        public.film_genre AS fg;
END
$$ LANGUAGE plpgsql;
```

# Using Functions



- Treated like a table
- SELECT
  - \*
  - Column list
  - JOIN

# Using Function Example

**SELECT**

*gl*.genre\_id,

*gl*.genrename

**FROM**

genre\_list() **AS** gl;

# Function Query Planner Parameters



- **COST** (100): estimated execution cost of the function
- **ROWS** (1,000): estimated number of rows from the function

# Function Query Planner Parameters



- **SET**: modify session configuration
  - search\_path is often set
- **PARALLEL {UNSAFE|SAFE}**
  - Unsafe (default) forces serial execution plan
  - Safe allows the planner to consider parallel plans

# Function Example

```
CREATE FUNCTION genre_list ()  
RETURNS TABLE (genre_id int, genre_name TEXT)  
COST 10000  
ROWS 100000  
AS $$  
BEGIN  
    RETURN QUERY  
    SELECT  
        fg.genre_id,  
        fg.name  
    FROM  
        public.film_genre AS fg;  
END  
$$ LANGUAGE plpgsql;
```

# Function Example

```
CREATE FUNCTION genre_list ()  
RETURNS TABLE (genre_id int, genre_name TEXT)  
SET search_path=bluebox,public  
AS $$  
BEGIN  
    RETURN QUERY  
    SELECT  
        fg.genre_id,  
        fg.name  
    FROM  
        film_genre AS fg;  
END  
$$ LANGUAGE plpgsql;
```

# Procedures



- Procedures are meant to perform actions, as opposed to functions, which return data
- Similar to Oracle procedures, or SQL Server functions

# Procedure Basics



- Use CREATE PROCEDURE command
- Provide a name
- Create a code block
- Define language used

# Procedure Example

```
CREATE OR REPLACE PROCEDURE new_genre (genrename TEXT)
AS $$

BEGIN
    INSERT INTO bluebox.film_genre
        (name) VALUES (genrename);

END
$$

LANGUAGE plpgsql;
```

# INOUT Procedure Example

```
CREATE OR REPLACE PROCEDURE
    genre_count(genreid int, INOUT genrecount int )
AS $$

BEGIN
    SELECT count(*)
    INTO genrecount
    FROM
        film AS f
    WHERE genreid = ANY(f.genre_ids) ;
END
$$
LANGUAGE plpgsql;
```

# Using Procedures



- CALL
- INOUT parameters

# Using Procedures Examples

```
CALL new_genre('Anime');

DO $$

DECLARE genre_counts int;
BEGIN
    CALL genre_count(genreid => 27, genrecount => genre_counts);
    RAISE NOTICE 'Films in genre: %', genre_counts;
END $$
```

# Subtransactions

- Internal transactions are committed separately
- Helpful for processing large batches of data over time
  - Checking progress



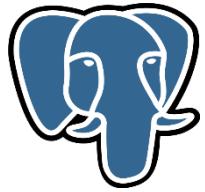
# Roles and Privileges

# The Basics

- **Roles** – The principle that can own objects and be granted permissions. E.g. Users & Groups
- **Privileges** – The set of permissions that a role can be granted to objects in a database or cluster
- **Ownership** – Every object in a cluster and database is owned by a role and carries a lot of power
- **Cluster** - The instance of PostgreSQL that's running, each containing one or more databases. A server can have multiple, independent clusters
- **Database** – The object we mostly care about. ;-) Also essential for database connections. Users connect to a database, not the cluster (specifically)



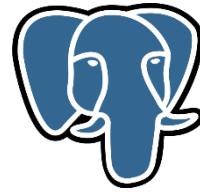
## Server/Host (Firewall, Ports)



**Cluster**

**Port: 5432**

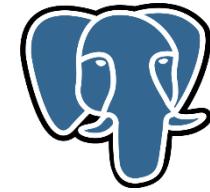
**pg\_hba.conf**



**Cluster**

**Port: 5433**

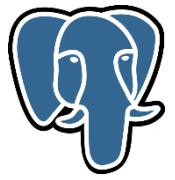
**pg\_hba.conf**



**Cluster**

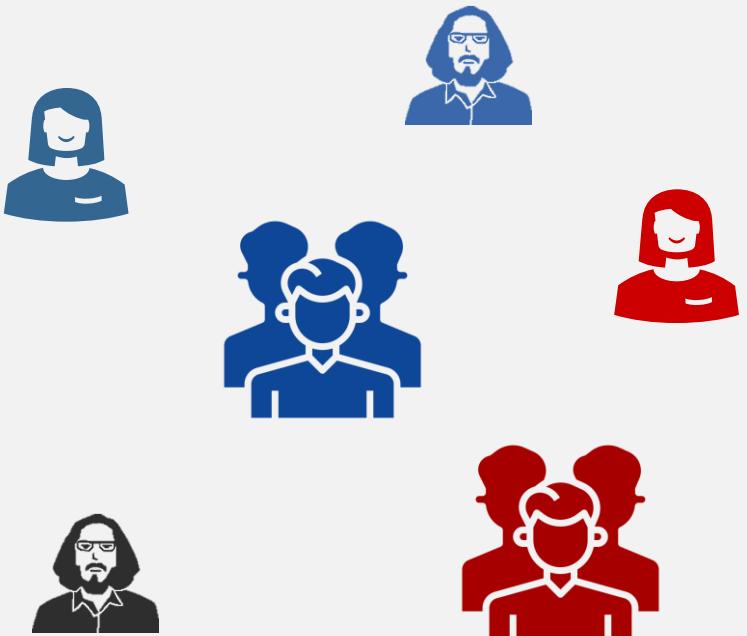
**Port: 5434**

**pg\_hba.conf**



# Cluster

## ROLES



## Databases



# `pg_hba.conf`

- First layer of authentication
- Similar to a firewall ruleset for PostgreSQL
- Cloud vendors largely manage this for you

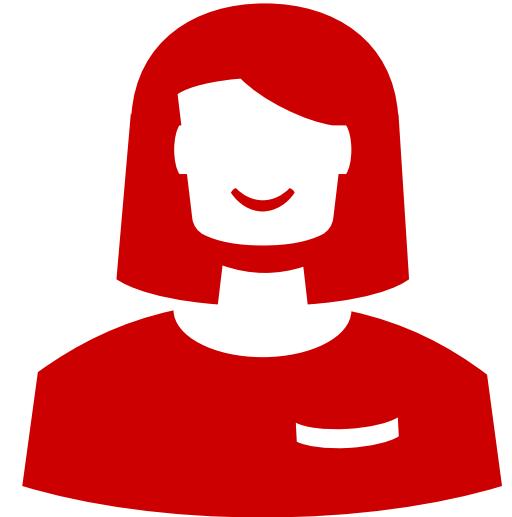
**Which hosts & roles, can connect to what databases,  
using what authentication method?**

```
# Allow any user on the local system to connect to any database with
# any database user name using Unix-domain sockets (the default for local
# connections).
#
# TYPE   DATABASE        USER        ADDRESS             METHOD
local   all            all
#
# The same using local loopback TCP/IP connections.
#
# TYPE   DATABASE        USER        ADDRESS             METHOD
host    all            all          127.0.0.1/32      trust
#
# Allow any user from host 192.168.12.10 to connect to database
# "postgres" if the user's password is correctly supplied.
#
# TYPE   DATABASE        USER        ADDRESS             METHOD
host    postgres        all          192.168.12.10/32  scram-sha-256
```

<https://www.postgresql.org/docs/current/auth-pg-hba-conf.html>

# Roles

- Own databases, schemas, and objects
  - Tables, Functions, Views, Etc.
- Have cluster-level privileges (attributes)
- Granted privileges to databases, schemas, and objects
- Can possibly grant privileges to other roles



# PostgreSQL 15 Attributes

LOGIN

SUPERUSER

CREATEROLE

CREATEDB

REPLICATION LOGIN

PASSWORD

INHERIT

BYPASSRLS

CONNECTION LIMIT

# Users and Groups

- Semantically the same as roles
- By Convention:
  - User = **LOGIN**
  - Group = **NOLOGIN**
- PostgreSQL 8.2+ **CREATE (USER | GROUP)** is an alias

```
CREATE USER user1 WITH PASSWORD 'abc123' INHERIT;
```

```
CREATE GROUP group1 WITH INHERIT;
```

```
CREATE ROLE user1 WITH LOGIN PASSWORD 'abc123' INHERIT;
```

# PostgreSQL Superuser

- 🦸 is created by default when the cluster is initialized
- Typically named `postgres` because the system process user initiates a `initdb`
- Bypasses all security checks except `LOGIN`
- Full privilege to do "anything"
- Treat superuser with care (like `root` on Linux)

**Most cloud providers do not  
provide superuser access**

# Superuser-like

- Create a role with the right level of control
- Recommend adding `CREATEROLE` and `CREATEDB`
- Allows user management and database ownership
- May still limit some actions (e.g. installing extensions limited to superuser)

# PostgreSQL 15+ Privileges

SELECT

INSERT

UPDATE

DELETE

TRUNCATE

REFERENCES

TRIGGER

CREATE

CONNECT

TEMPORARY

EXECUTE

USAGE

SET

ALTER SYSTEM

# PUBLIC Role

- All roles are granted implicit membership to PUBLIC
- The public role cannot be deleted
- Granted CONNECT, USAGE, TEMPORARY, and EXECUTE by default
- <=PG14: CREATE on the public schema by default
- >=PG15: No CREATE on public schema by default

# Granting Privileges

```
-- grant the ability to create a schema  
GRANT CREATE ON DATABASE app_db TO admin1;
```

```
-- see and create objects in schema  
GRANT USAGE,CREATE IN SCHEMA demo_app TO dev1;
```

```
-- allow some roles only some privileges  
GRANT SELECT,INSERT,UPDATE  
ON ALL TABLES IN SCHEMA demo_app TO jr_dev;
```

# Granting Privileges

- Remember, explicit grants only effect existing database objects!

```
-- This will only grant to existing objects  
GRANT ALL TO ALL TABLES IN SCHEMA public TO dev1;
```

# More Detail on GRANT and REVOKE

**What the privileges mean:**

<https://www.postgresql.org/docs/current/ddl-priv.html>

**How to GRANT privileges:**

<https://www.postgresql.org/docs/current/sql-grant.html>

**How to REVOKE privileges:**

<https://www.postgresql.org/docs/current/sql-revoke.html>

**DEMO!**



# PASS Session

## For Your Eyes Only: Roles, Privileges, and Security in PostgreSQL

Wednesday, Nov. 6  
9:45-10:45  
Room 445-446



<https://www.youtube.com/watch?v=3bCWLe-dU58>

# MVCC and Vacuum

# **'MVCC Unmasked' by Bruce Momjian**

Portions of this content were sourced from Bruce's seminal talk, "**MVCC Unmasked**", under the Creative Commons Attribution License.

<https://momjian.us/main/writings/pgsql/mvcc.pdf>

PostgreSQL documentation uses the word tuple for row.

In technical terms, a tuple is an instance of the row schema.

Multi-Version Concurrency Control (MVCC) is what allows PostgreSQL (and other databases) to support highly concurrent workloads with the goal that:

*"readers never block writers,  
and writers never block readers"*

# Multi-Version Concurrency Control

For the default **Read Committed Isolation Level**:

- Each query only sees transactions that were committed before it started (ID less than current transaction)
- At the beginning of a query, PostgreSQL records:
  - The transaction ID (xid) of the current command
  - All transaction IDs that are in process
- Multi-statement transactions also record any of their own completed statements from earlier in transaction

# Assigning XIDs

In PostgreSQL, transaction IDs are assigned when the first command/statement is run within a transaction, not when **BEGIN** is executed.

# INSERT

- xmin = current transaction ID
- xmax = zero (invalid)

xmin	xmax	User Data
48	0	Hello!

## DELETE

- xmin is unchanged
- xmax = current transaction ID

Dead Tuple

xmin	xmax	User Data
48	78	Hello!

# UPDATE

Original tuple:

- xmin is unchanged
- xmax = current transaction ID

New tuple:

- xmin = current transaction ID
- xmax = zero

Dead Tuple

xmin	xmax	User Data
48	78	Hello!
78	0	Hello World!

## UPDATE (with the same value)

Original tuple:

- xmin is unchanged
- xmax = current transaction ID

New tuple:

- xmin = current transaction ID
- xmax = zero

Dead Tuple

xmin	xmax	User Data
48	78	Hello!
78	0	Hello!

# Visible Tuples

A tuple is considered "visible" when creation XID (**xmin**):

- is a committed transaction
- is less than the XID stored at the start of the current query
- was not in-process at the start of the current query

Visible tuples must *also* have an expired XID (**xmax**) that:

- is blank *or* aborted *or*
- is greater than the transaction counter stored at query start *or*
- was in-process at query start

# Tuple Visibility and Snapshots

## Created rows only

xmin	xmax	User Data
10		Hello!
50		Velkommen
100		willkommen

Visible  
Not Visible  
Not Visible

## Snapshot

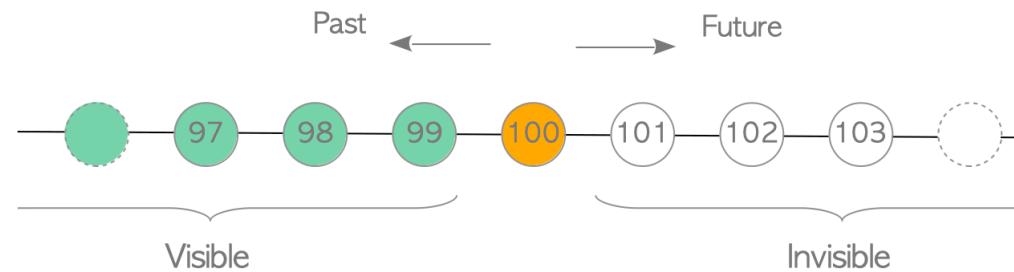
- Highest committed XID = 90
- Open transactions: 20,50,100
- Current XID = 110

## Updated rows

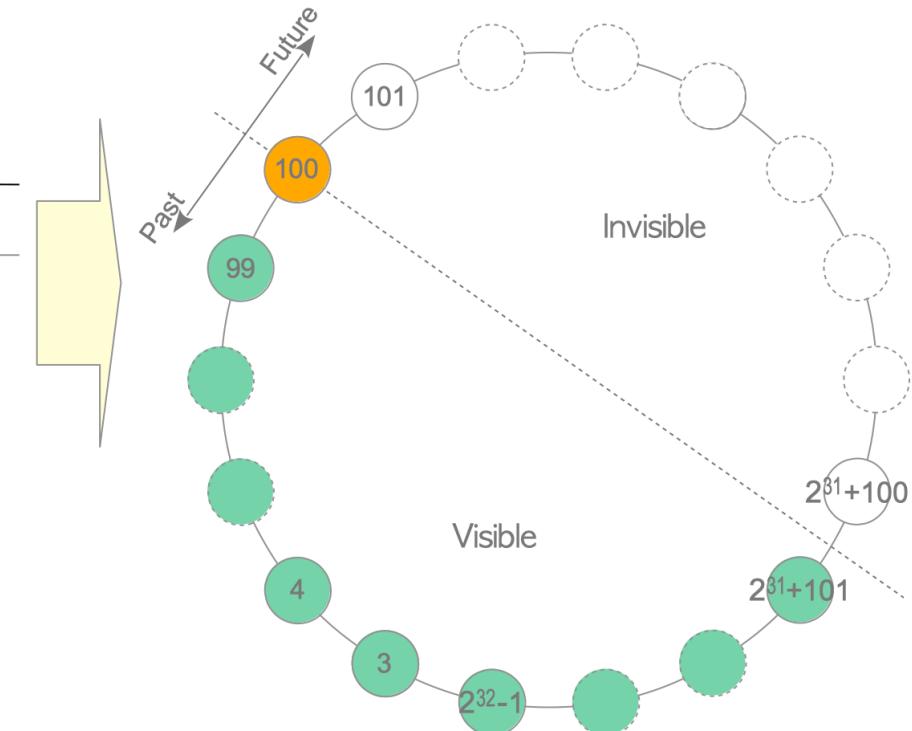
xmin	xmax	User Data
10	70	Hello!
40	50	Velkommen
48	100	willkommen

Not Visible  
Visible  
Visible

a) Transaction identifiers



b) Transaction identifiers space as a circular



## XIDs are stored as INT (32-bit)

- Compared using modulo- $2^{32}$
- 2.1 billion XIDs in the "past"
- 2.1 billion XIDs in the "future"

Image Credit: <https://www.interdb.jp/pg/pgsql05.html>

# PostgreSQL Heap Page Layout

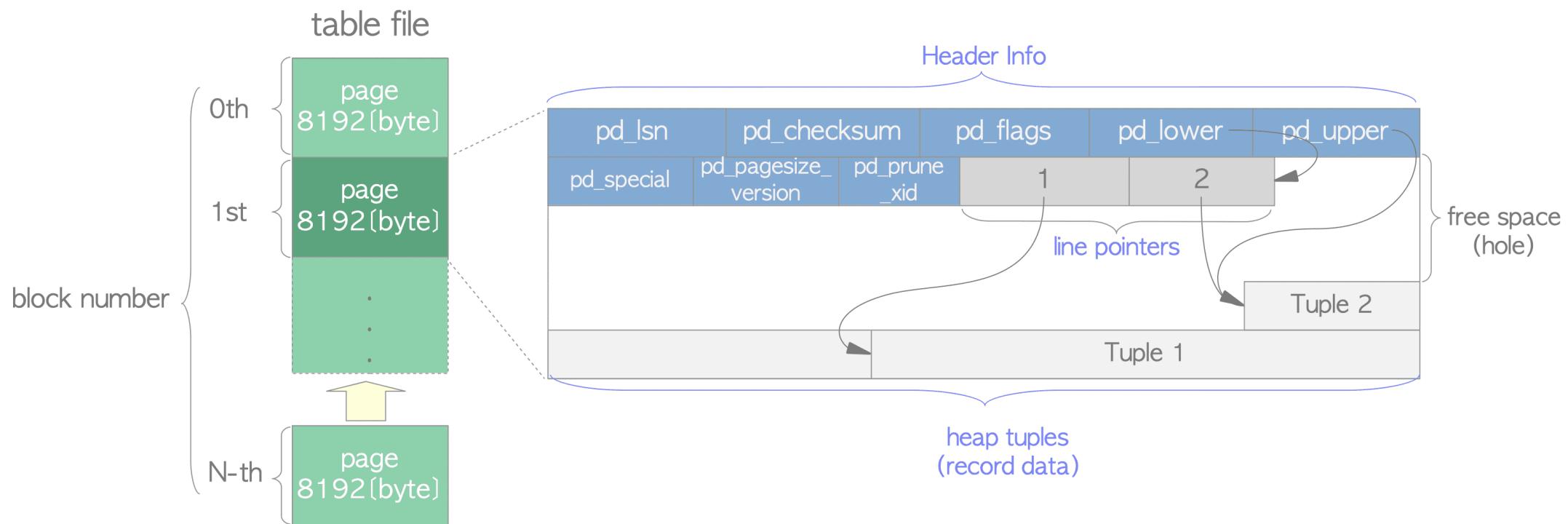
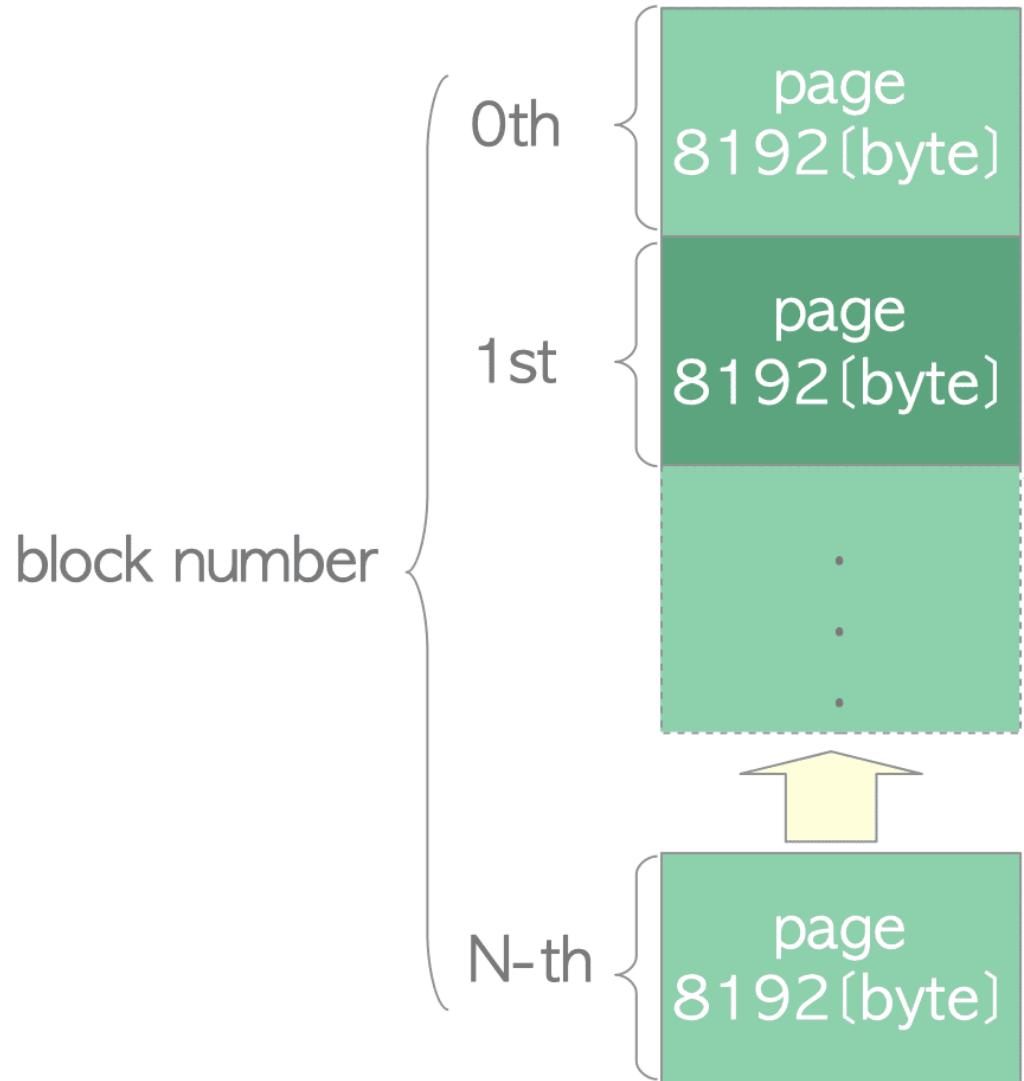
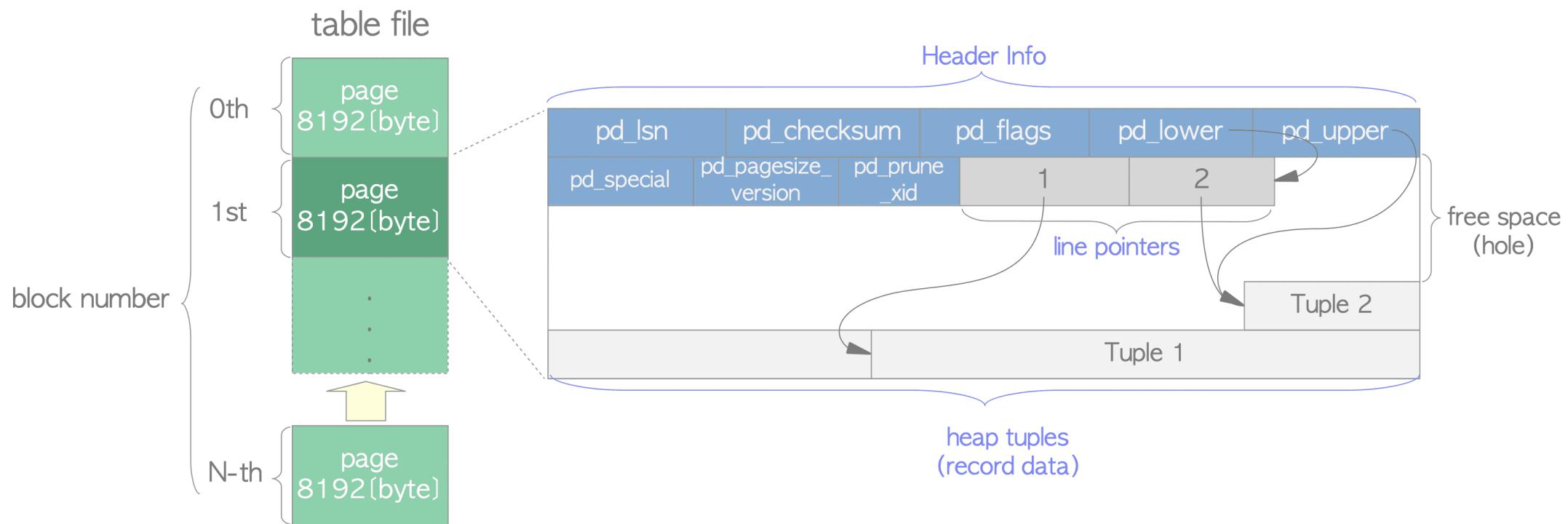


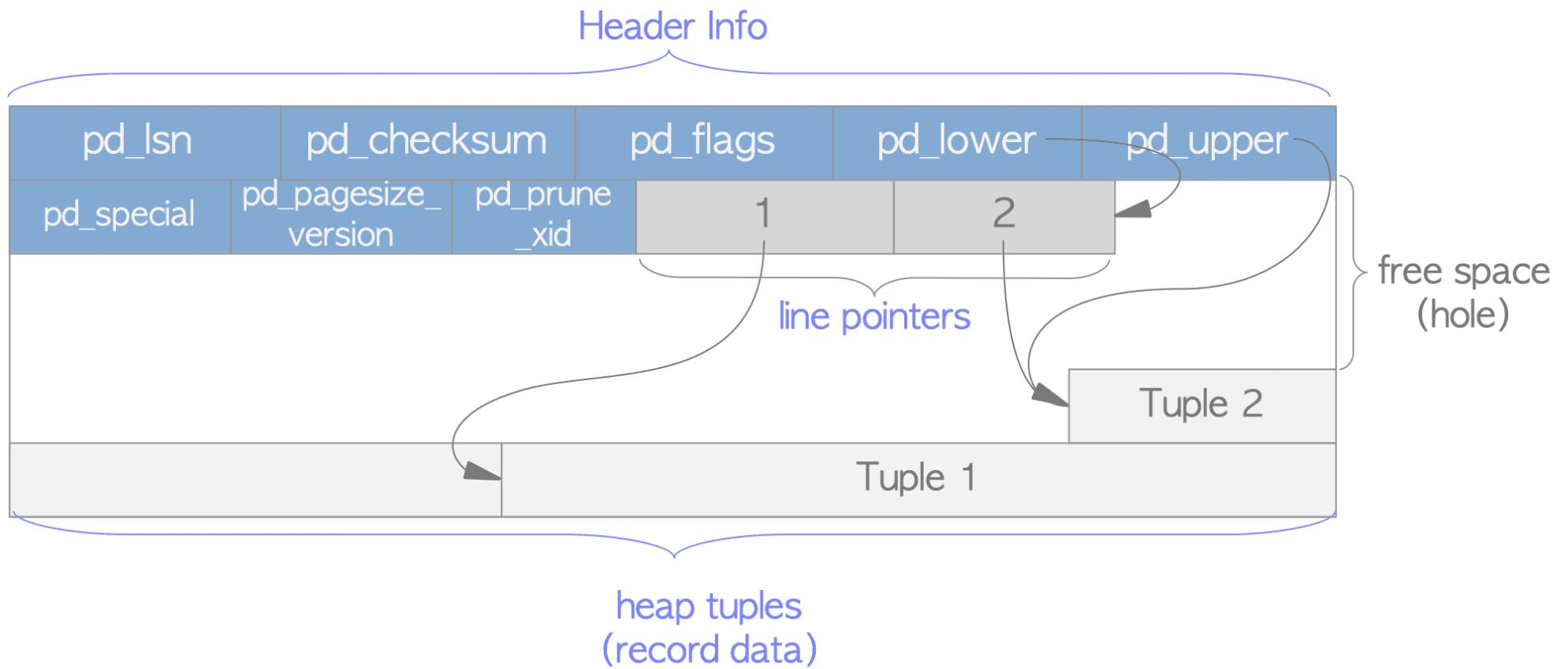
Image credit: <https://www.interdb.jp/pg/pgsql01.html>

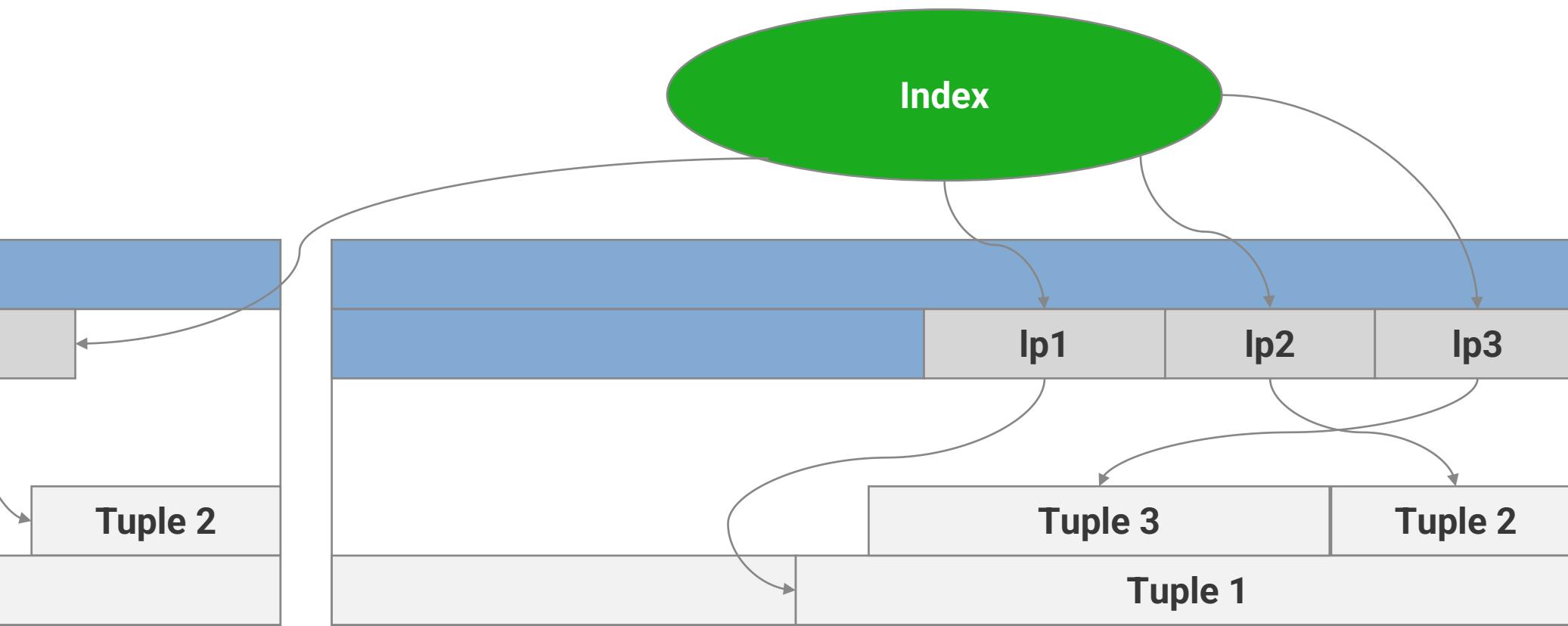
# table file

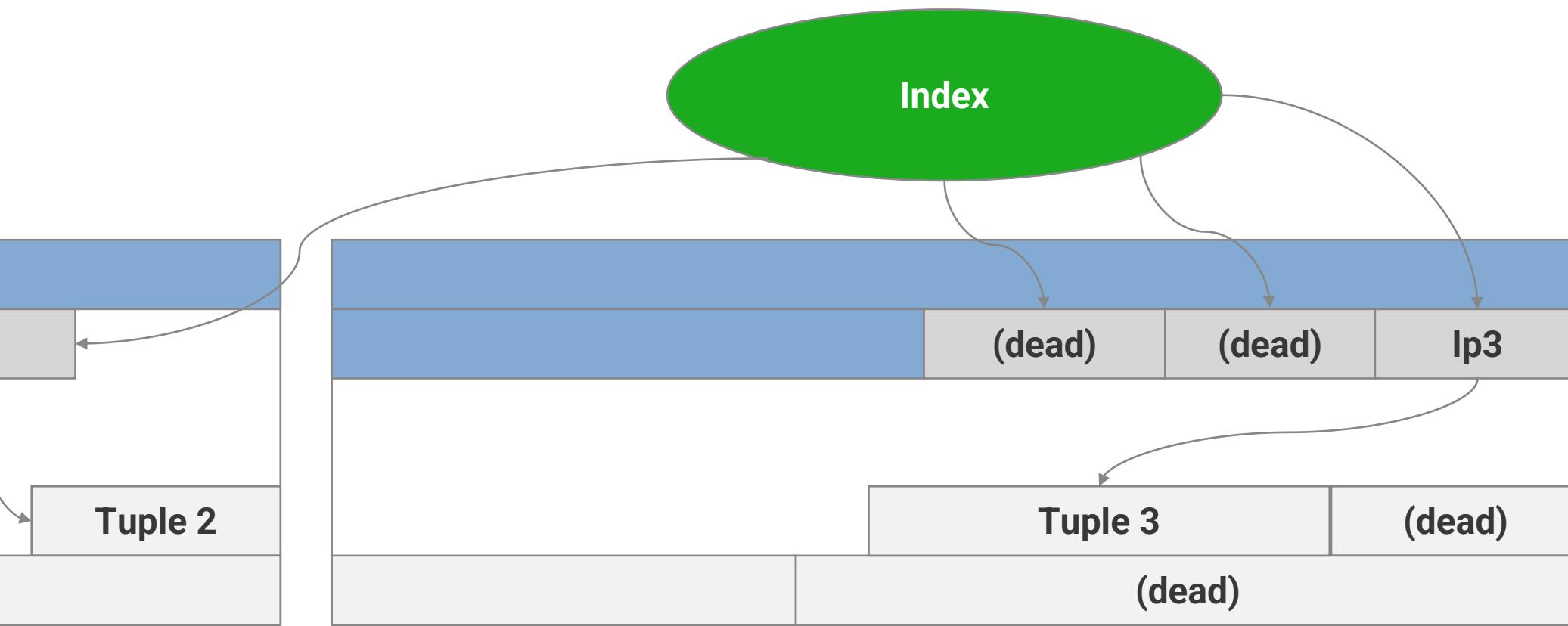


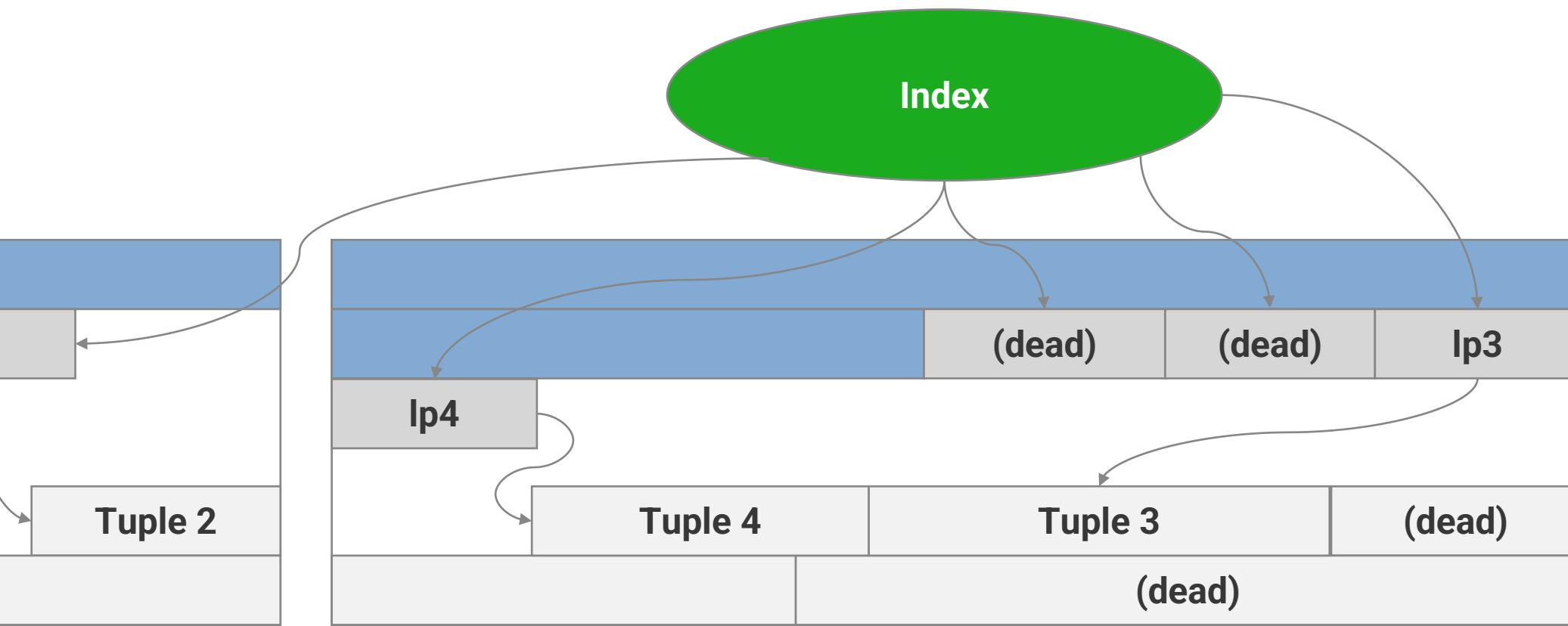
# PostgreSQL Heap Page Layout

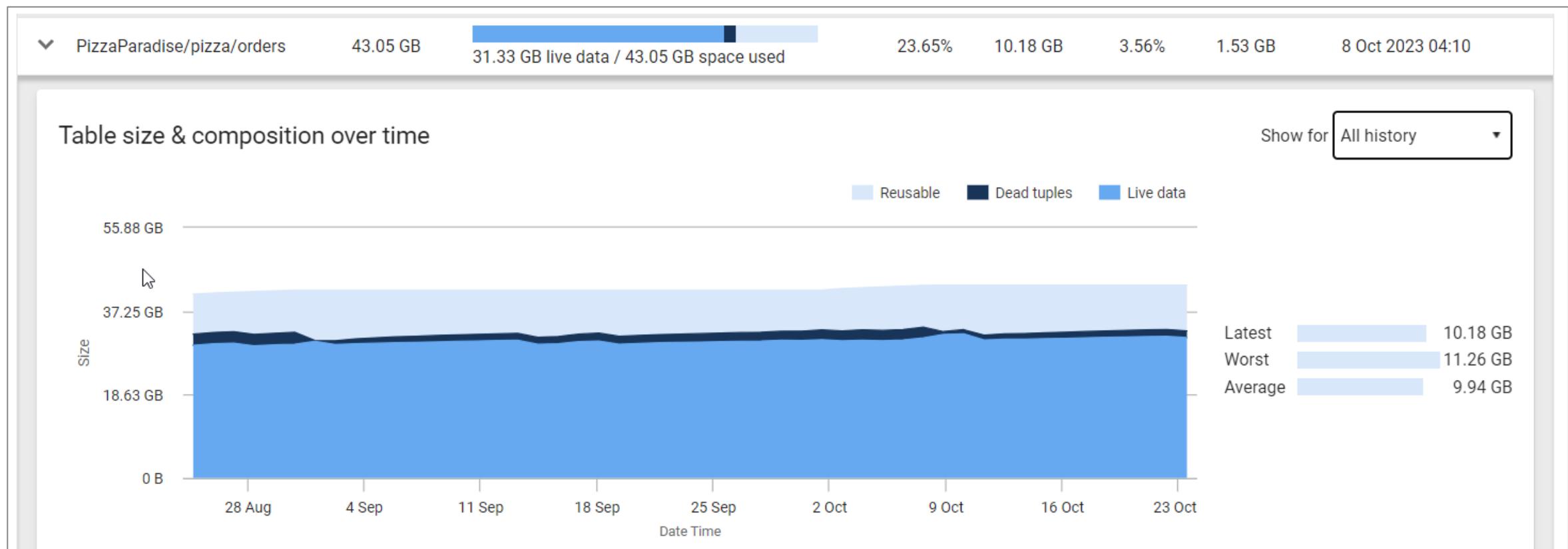




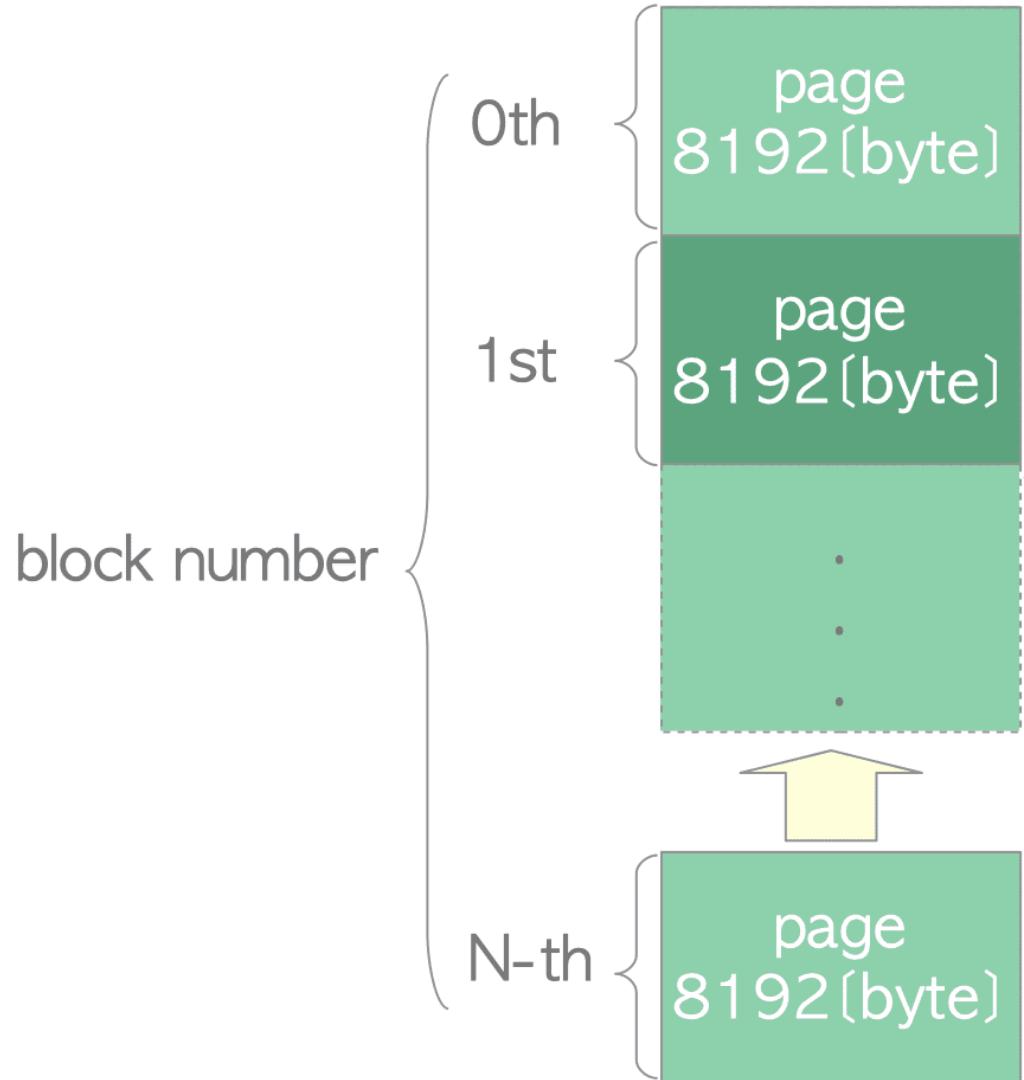






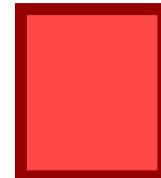


# table file

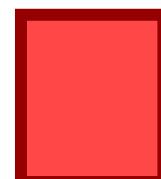
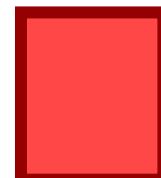
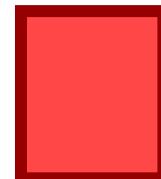
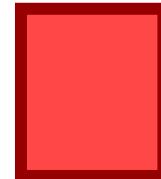
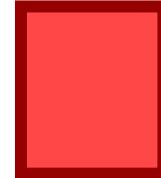


xmin	xmax	ID	User data...

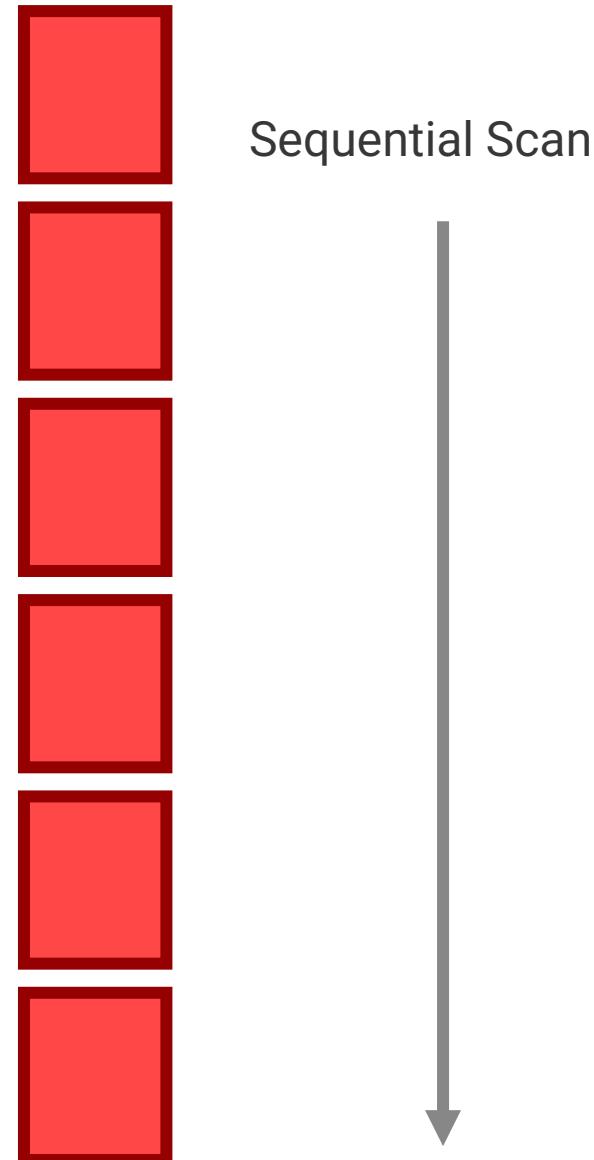
SELECT a, b, c FROM t;



Sequential Scan

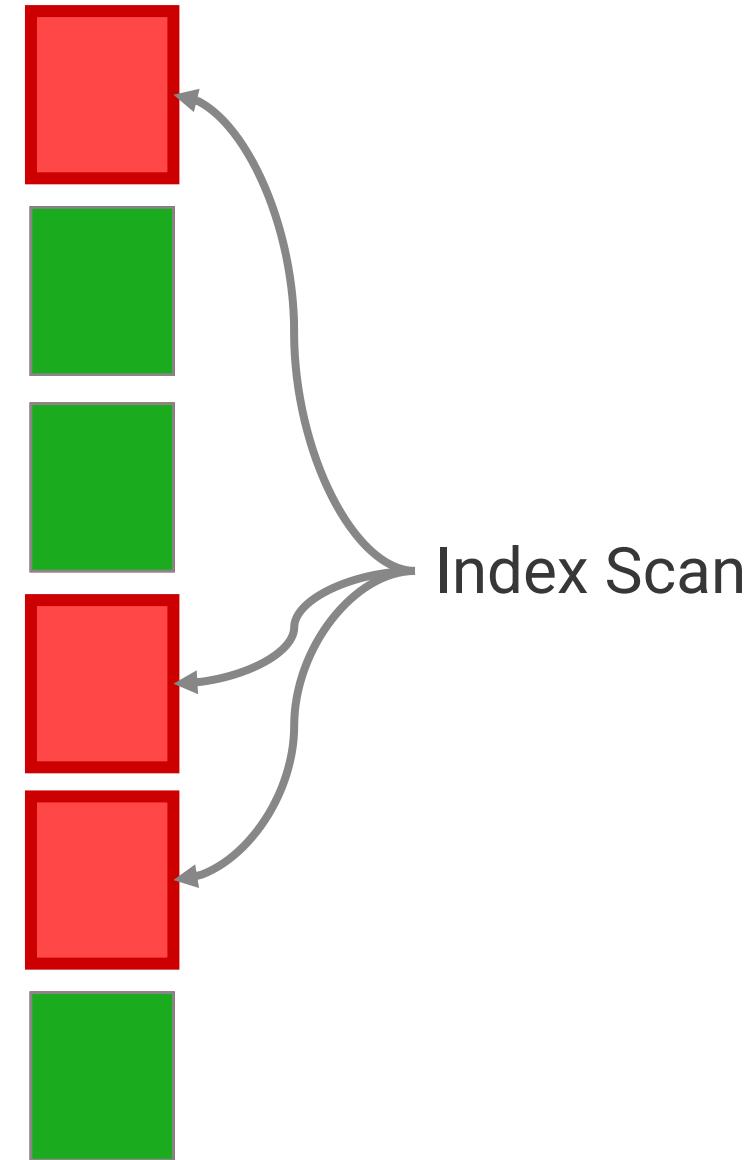


```
SELECT a, b, c FROM t  
WHERE ID=10;
```



xmin	xmax	ID	User data...

SELECT a, b, c FROM t  
WHERE ID=10;

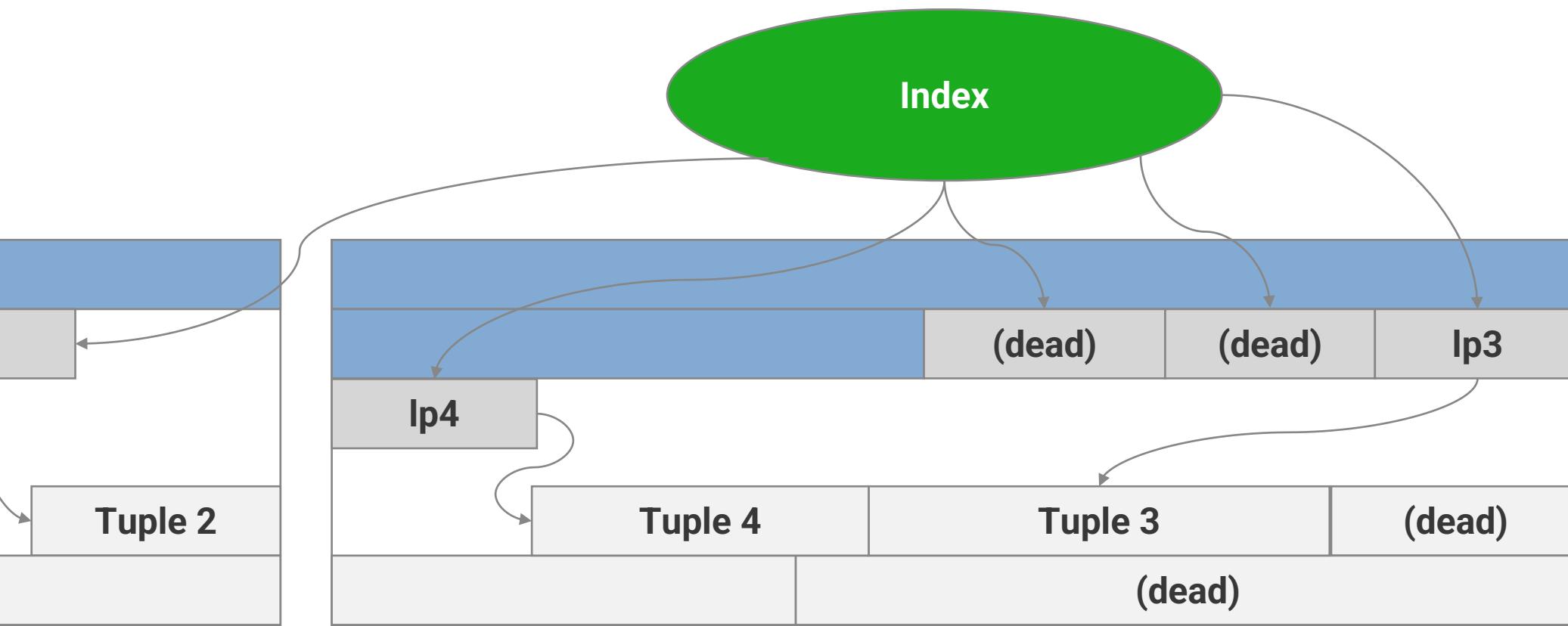


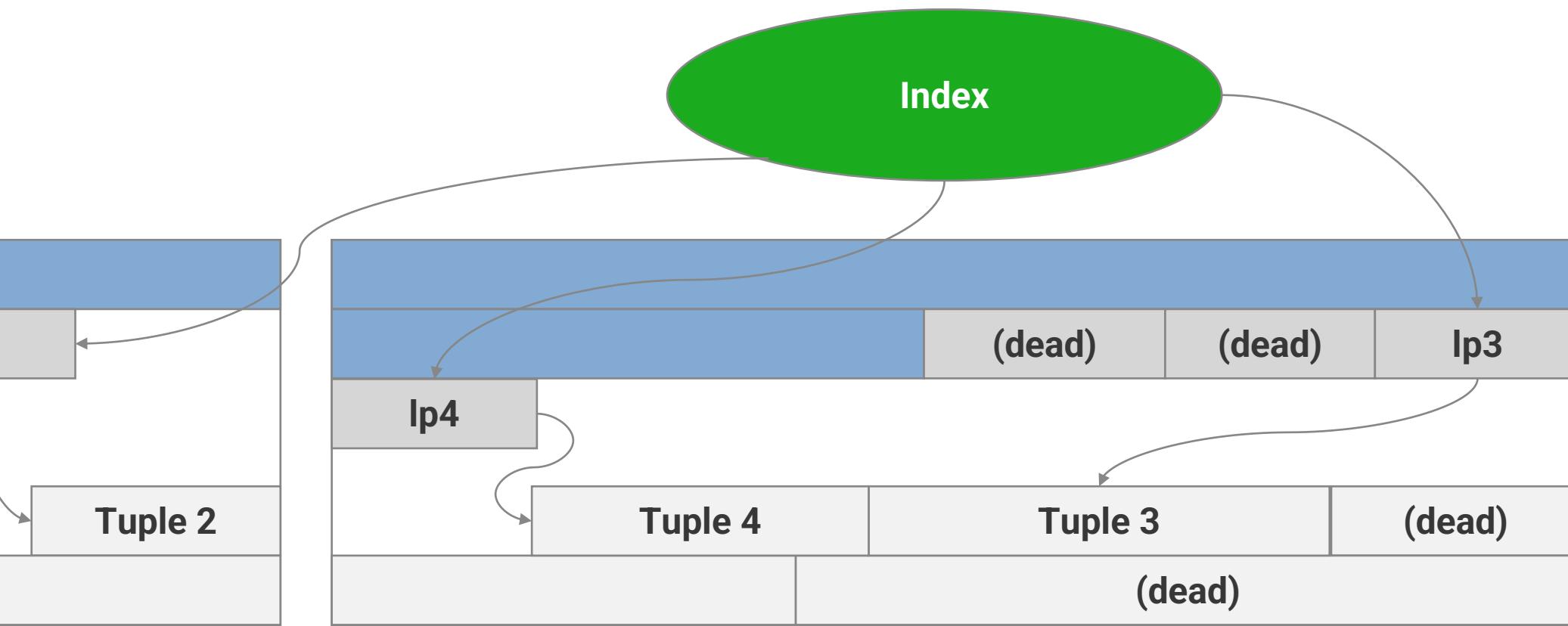


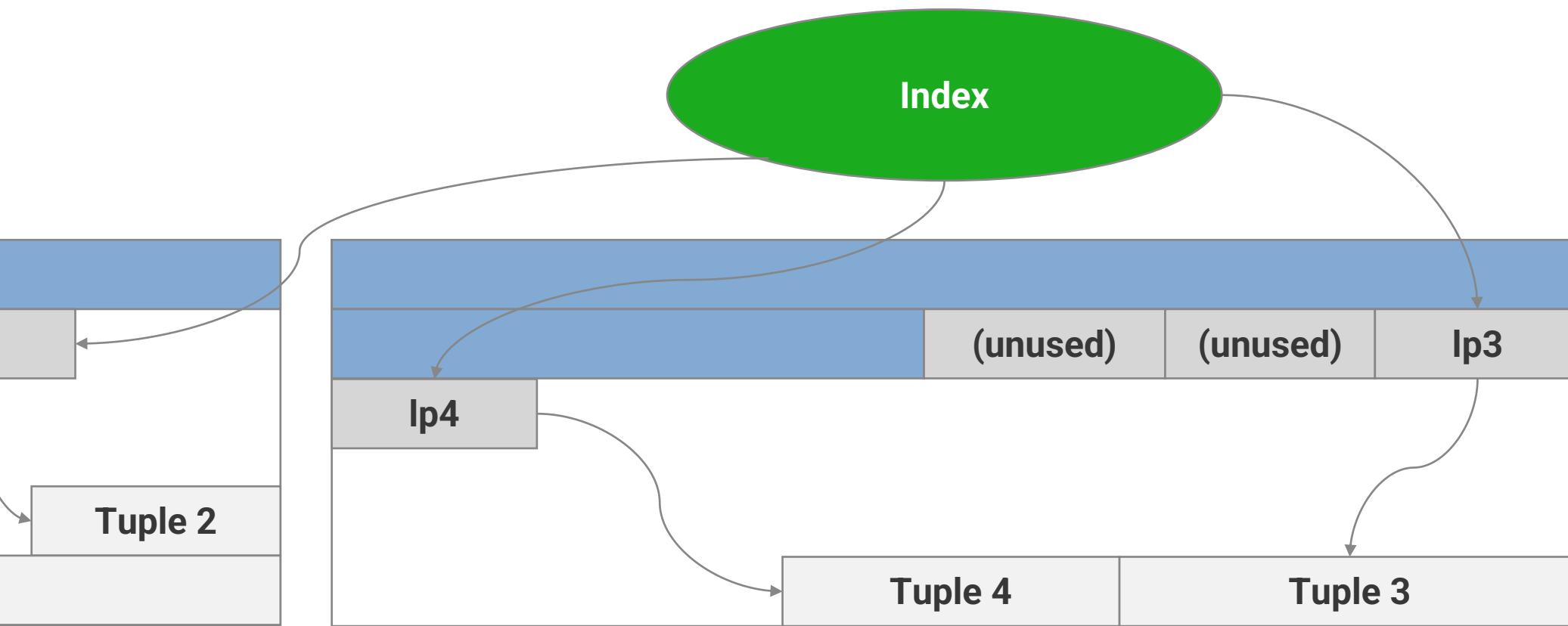
# Vacuum!

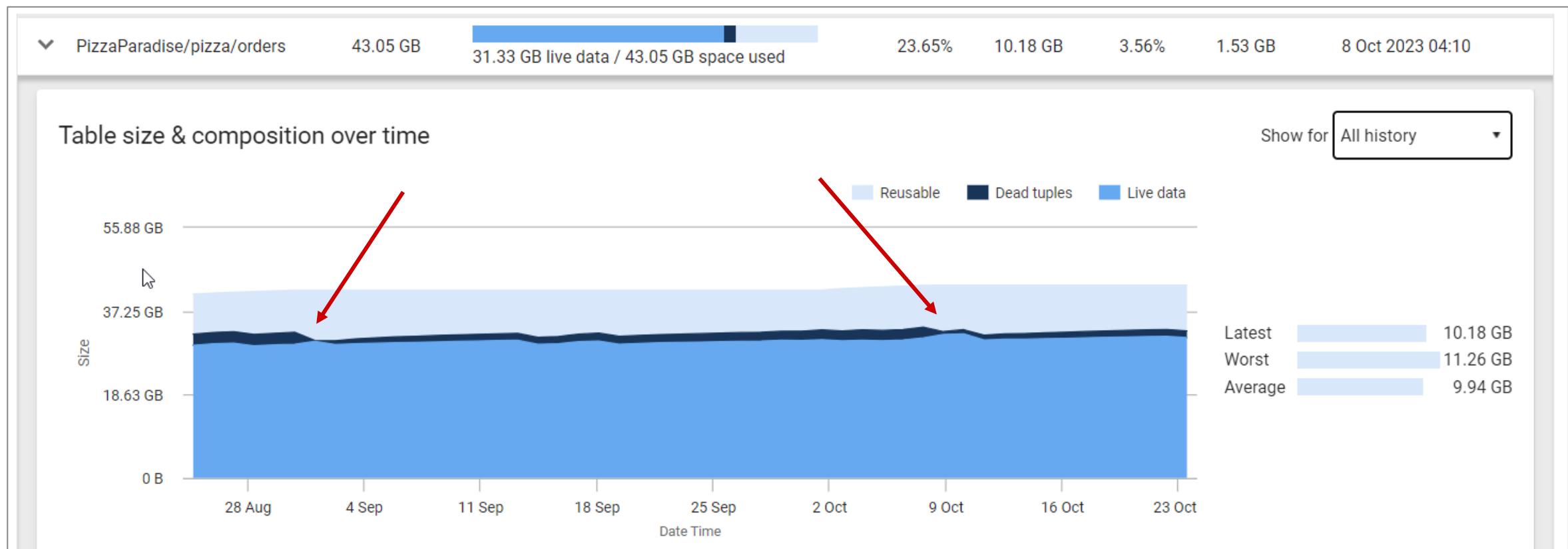
- Serves two main purposes:
  - Free up space on pages used by dead tuples
  - Mark rows as "frozen" if they have a XID sufficiently older than the most recent XID, allowing the XID to be reused again in the future
  - Additionally, if all rows on a page are frozen, the page can be marked frozen to speed up scans

Free space used by  
dead tuples

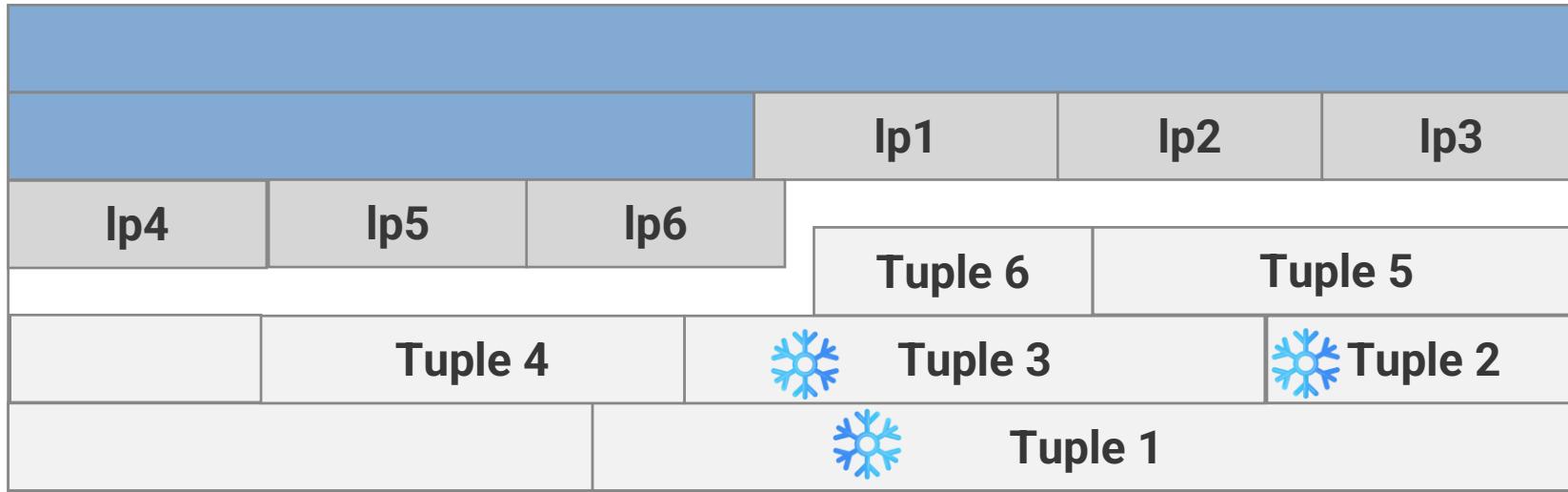


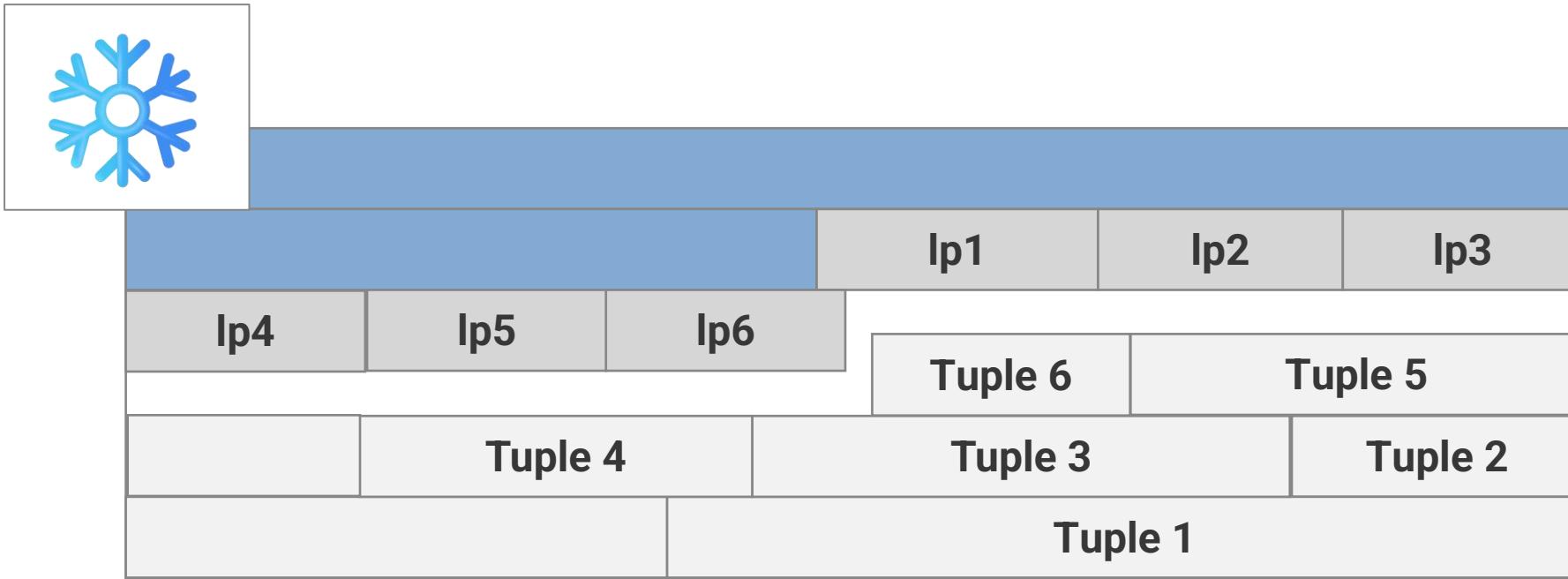






# Freeze rows and pages





# VACUUM Automation

- Vacuuming is so vital that there is an automated process called **autovacuum**
- Triggered based on a percentage of dead tuples/modified data (generally)
- Individual DELETE or UPDATE heavy tables can be tuned directly

# Individual Table Auto(vacuum/analyze)

- Autovacuum triggers at 20% data modification
- Autoanalyze triggers at 10% data modification
- Analogous to TF-2731, but per table!

# Global Autovacuum and Autoanalyze Settings

"Scale Factor" = Percentage of rows modified before process runs

"Threshold" = Added number of rows to add to the percent calculation.  
(this prevents very small tables from taking resources)

ანუ ფორმული შედგება საჭირო გარემოს გასაკვირვებლად . , , .

ანუ ფორმული შედგება საჭირო გასაკვირვებლად .

ანუ ფორმული ანალიზის გარემოს გასაკვირვებლად . , , .

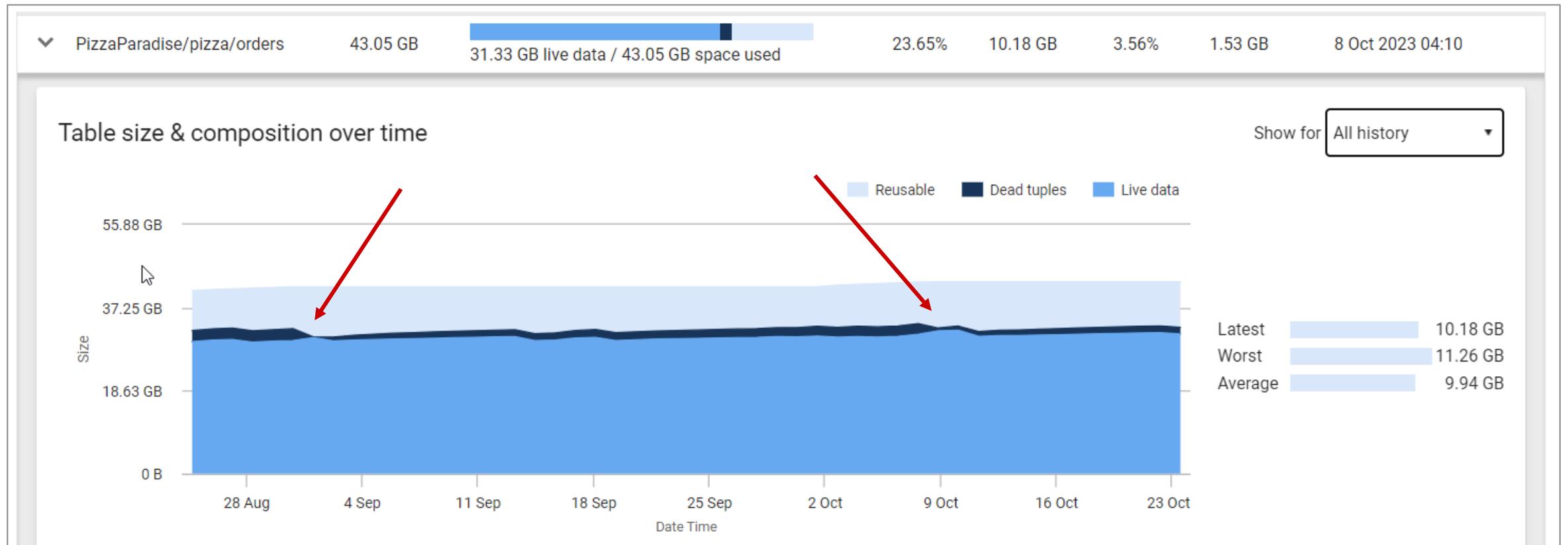
ანუ ფორმული ანალიზის გასაკვირვებლად .

**Analyze at 5% modification:**

```
ALTER TABLE t SET (autovacuum_vacuum_scale_factor=0.05)
```

**Analyze threshold above percentage:**

```
ALTER TABLE t SET (autovacuum_vacuum_threshold=500)
```



**'If you think autovacuum/analyze is  
running too frequently, it's probably  
not running frequently enough'**

PASS Session

# Importance of PostgreSQL Vacuum Tuning to Optimize Database Performance



Wednesday, Nov. 6



2:00-3:00



Room 447



# Standard Configuration and Maintenance

Over 300 configuration  
options

# Configuration

- *Every non-serverless Postgres should be tuned*
- Edit `postgresql.conf`

## Helpful Tuning Resources

[https://wiki.postgresql.org/wiki/Tuning\\_Your\\_PostgreSQL\\_Server](https://wiki.postgresql.org/wiki/Tuning_Your_PostgreSQL_Server)

<https://www.youtube.com/watch?v=IFIpm73qtk>

<https://postgresqlco.nf/>

# Configuration – The Big 4

## **shared\_buffers:**

- data cache
- at least 25% of total RAM

## **work\_mem:**

- max memory for each plan operation
- JOINS, GROUP BY, ORDER BY can all have workers

# shared\_buffers

- Configurable memory reserved for data cache
- Start at 25% of total server memory
- Keep a watch on Cache Hit Ratio
  - <90% consistently, increase shared\_buffers or increase instance resources

# Cache Hit Ratio

SELECT

```
    sum(heap_blks_read) as heap_read,  
    sum(heap_blks_hit) as heap_hit,  
    sum(heap_blks_hit) / (sum(heap_blks_hit) +  
    sum(heap_blks_read)) as ratio
```

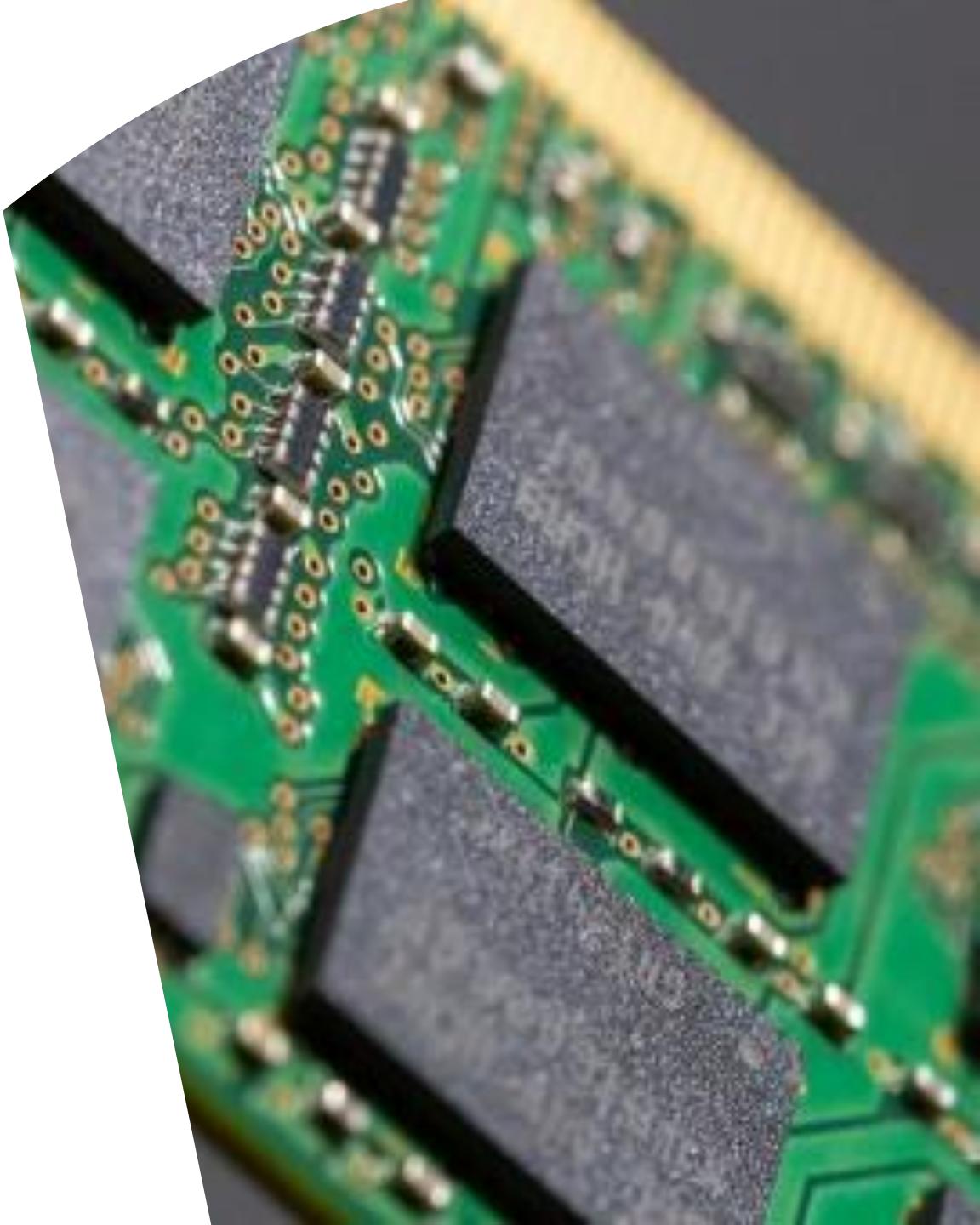
FROM

```
pg_statio_user_tables;
```

# work\_mem

- Single most important setting for complex workloads\*
- Memory used by query operations like sorting and hashing

\*It Depends...



# `work_mem`

- Postgres doesn't pre-allocate memory before running a query (i.e. memory grant)
- When you see "external disk" operations in query plans, `work_mem` likely needs to be increased
- 4MB default, typically tune higher (8/16MB)
- Set on individual query sessions when necessary

# Configuration – The big 4

## **maintenance\_work\_mem:**

- memory for background tasks like VACUUM and ANALYZE
- Index maintenance tasks

## **max\_connections:**

- Defaults to 100
- Many providers start at 200
- Total memory usage equals:

**`max_connections * work_mem * hash/sort operation`**

# **maintenance\_work\_mem**

- **maintenance\_work\_mem** is the amount of RAM that a maintenance process like autovacuum can use to process pages of data. If it is not sufficient, vacuuming (and other essential maintenance processes) will take more time and fall behind.
- If autovacuum isn't completing its work regularly, check this setting!

**Default setting** = 64MB

# **autovacuum\_work\_mem**

- Amount of memory used by each autovacuum process
- Correlated to **autovacuum\_max\_workers**
- Tune if **maintenance\_work\_mem** is increased dramatically ( $>=1\text{GB}$ )

**Default setting** = **maintenance\_work\_mem**

# Configuration [other]

## **wal\_compression:**

- Off (default)
- Default uses pglz
- PG15 can use lz4 or zstd

## **autovacuum:**

- On (default)
- Do not turn off
- Might have to tune individual tables

# Configuration [other]

## **autoanalyze:**

- On (default)
- Do not turn off
- Might have to tune individual tables

## **default\_toast\_compression:**

- pglz (default)
- Consider lz4 instead

# Configuration (other)

## `random_page_cost`:

- "scale" of cost for random page access
- Defaults to 4.0 – spinning rust
  - Increases likelihood of table scans
- Change to 1.0 in most environments and monitor

# SHOW & SET

- **SHOW** will display the current value
- **SET** allows settings to be modified (when possible)
- Examples:
  - **SHOW ALL** – page all settings
  - **SHOW max\_connections** - configured connection limit

# Other Maintenance Tasks

# Reindexing



- Indexes can have bloat, too
- Reindexing consolidates and cleans index pages
- Requires access exclusive lock, unless CONCURRENTLY is used

# Reindexing



- Maintains FILLFACTOR when reindexing
- Essential upkeep for some specific index types like BRIN

# FILLFACTOR

- The target amount of free space left on each data page
- Defaults:
  - Tables = 100%
  - Indexes = 90%
- Very important consideration for update-heavy tables



# FILLFACTOR

## For indexes

- Higher value when indexed columns will not change (i.e. timestamps or counters)
- Lower value when data inserts will likely produce many page splits

## For tables

- Higher value for insert-only workloads
- Lower value for update heavy tables
  - Remember, every update produces a new row version

# Heap-only Tuples [HOT]

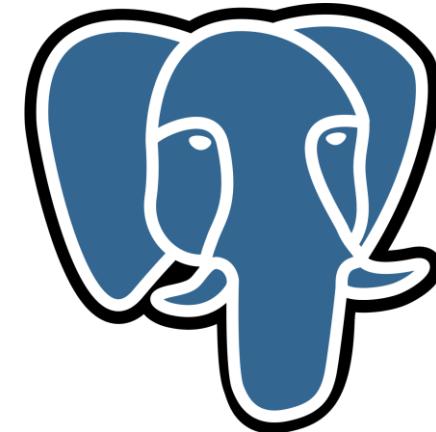
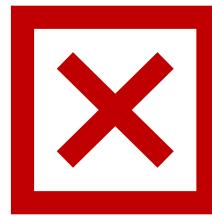
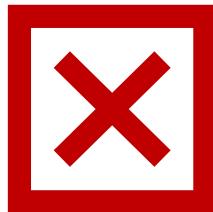
- HOT updates attempt to store new row version on the same data page
- Index pointer does not have to change
- Updated value cannot be referenced by an index
- Increase likelihood by decreasing FILLFACTOR
- Benefits update-heavy tables

```
SELECT schemaname, relname,  
       n_tup_upd, n_tup_hot_upd  
FROM pg_stat_user_tables;
```

	schemaname	relname	n_tup_upd	n_tup_hot_upd
1	bluebox	rental	4,489	0
2	bluebox	inventory	4,461	442

# Extensions

# Customizable with Extensions



# What are extensions?

- Non-core code that can be installed into a database
- Add new functions or features
- Modify PostgreSQL runtime based on hooks within the code
  - Query planner
  - Data retrieval
  - Storage engine
  - More...

# What are extensions?

- Written in any supported language, including SQL
- Many popular extensions are written in C to gain the best performance
- The pgrx framework is a popular Rust-based toolchain for building extensions

# Problems that extensions can solve

- Having a repeatable set of functions available in a database for DBAs within your company
- Providing access to server data that hasn't been exposed yet
- Modifying data storage
- Influencing the query planner

If you wish PostgreSQL could do something, it's probably possible with an extension... within reason



TimescaleDB



Citus



PostGIS



# What are extensions?

- Extensions must be installed **per-database**
  - There are no "server-level" extensions
- PostgreSQL ships with 'contrib' extensions
- Others must be installed on the server
  - Manually (source, RPM/YUM packages)
  - PostgreSQL Extension Network (PGXN)
- Trusted vs. Untrusted

# Extensions

See available extensions

```
SELECT * FROM pg_available_extensions ORDER BY name;
```

Install

```
CREATE EXTENSION pg_stat_statements;
```

Update

```
ALTER EXTENSION pg_stat_statements UPDATE;
```

Uninstall

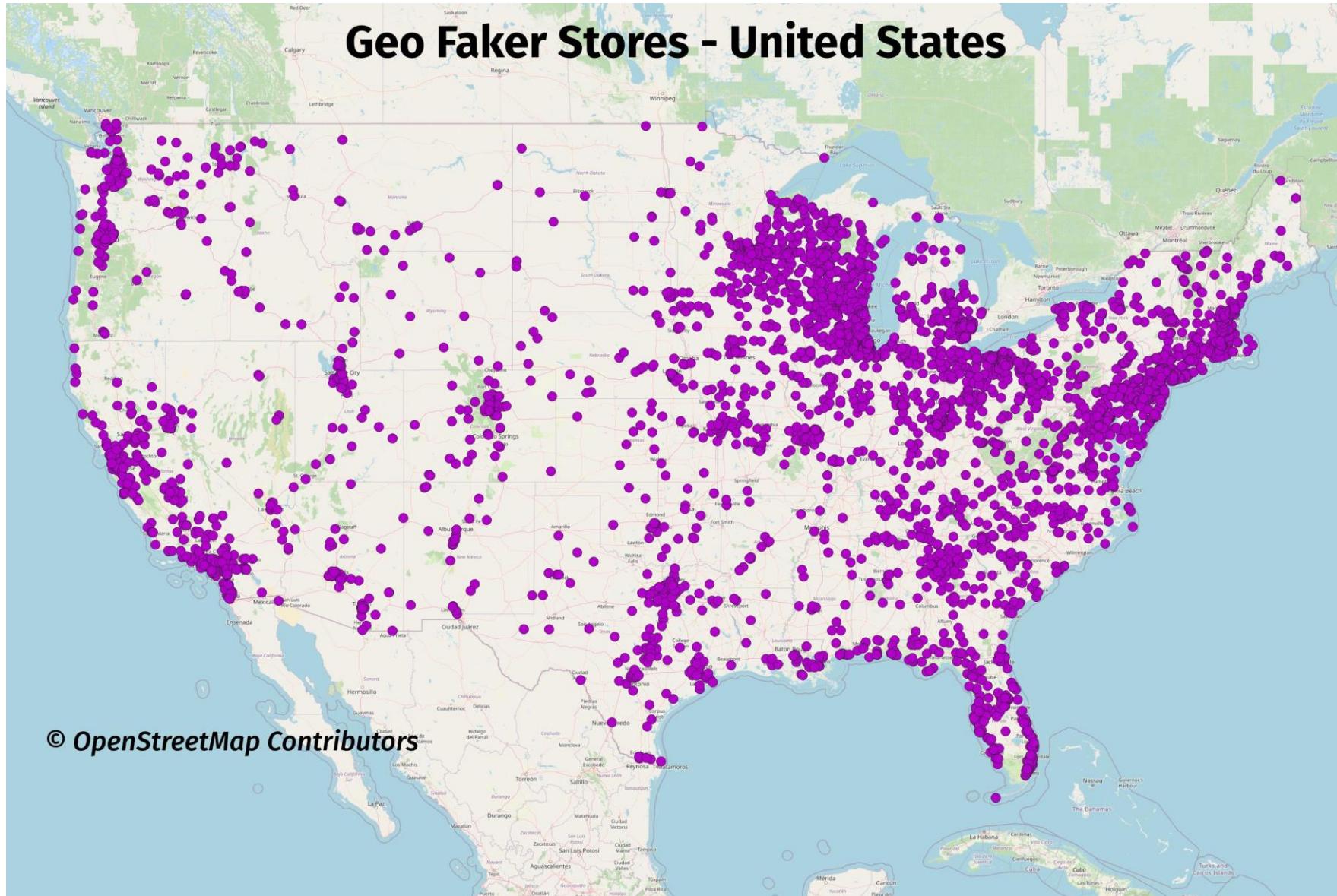
```
DROP EXTENSION pg_stat_statements;
```

Support varies among cloud  
providers

# 10 Extensions to know about

# PostGIS

- Industry-standard, open-source, geographic database
- Adds geographic:
  - Types: point, line, polygon
  - Indexes: GIN, GIST, SPGIST, BRIN
  - Functions: hundreds of "ST\_..." functions



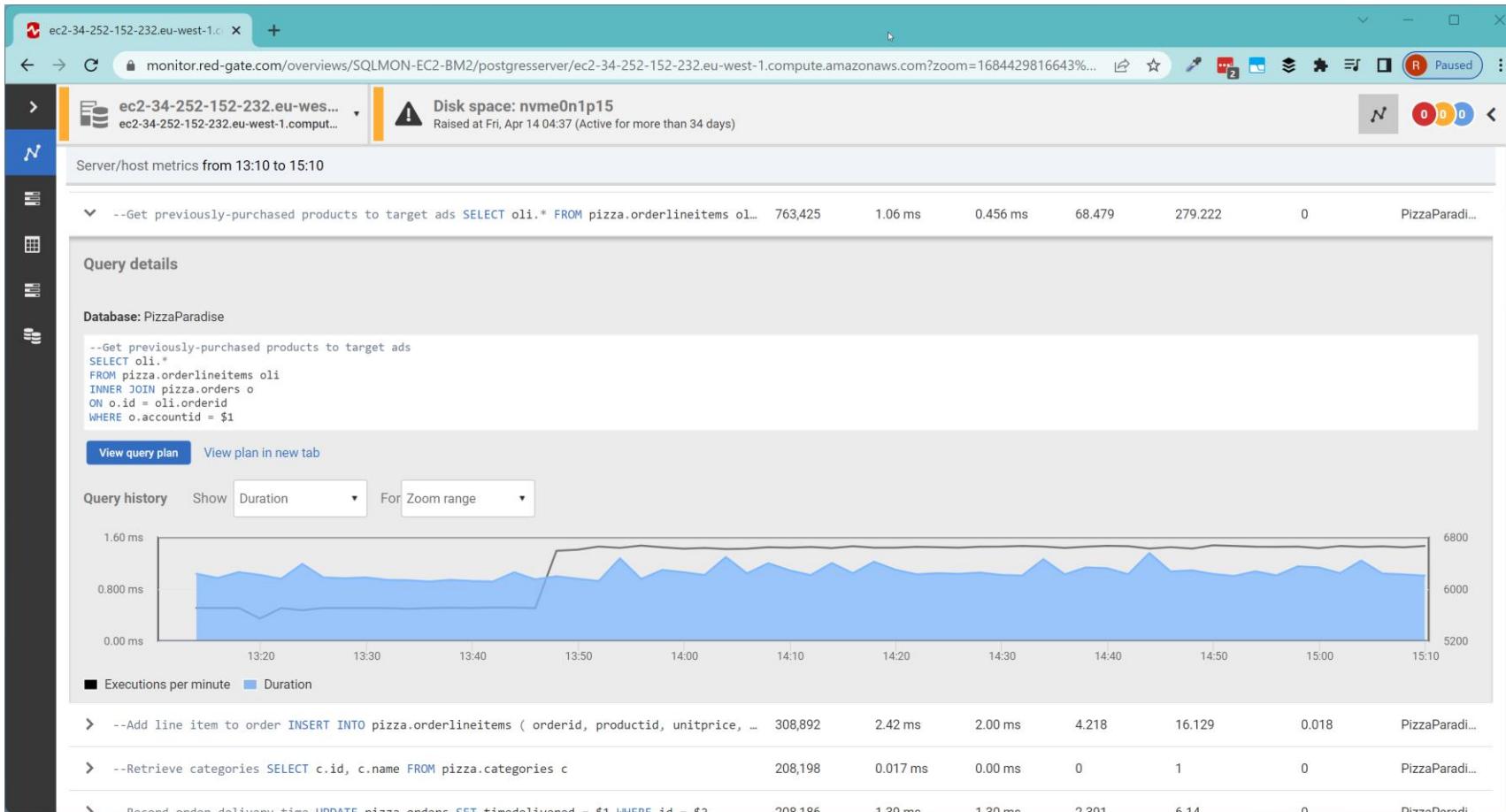
# postgres\_fdw

- Foreign Data Wrapper (FDW) is a standards-compliant connection infrastructure
- Initiates a connection to an external PostgreSQL database
- Necessary for cross-database *and* cross-instance queries
  - There is no fully-qualified notation to access other databases

# pg\_stat\_statements

- One of many statistics views within PostgreSQL
- Data is stored by PostgreSQL, view provides access to the data
- Matches counts for parameterized versions of queries
- Cumulative only, requiring tooling to store snapshots over time
- Installed by default on all major cloud providers

# SQL Monitor



# pg\_cron

- PostgreSQL does not have a built-in job scheduler
- pg\_cron is a popular extension for scheduling and running jobs on a schedule in the background
- Other extensions will often recommend pg\_cron for background work with that extension (e.g. pg\_partman)
- Available in all major cloud providers

# pg\_trgm

- Provide a new Trigram index type
- Splits text into groups of three characters
- Indexes for faster searching and relevance
- Indexes are large and slow, by comparison, so use appropriately
- Use for smaller text fields, like names or people or items

# hypopg

- Create an index that doesn't get materialized
- EXPLAIN the query to see if the planner would use the new index
- Supports many different types of index variations
- Extremely useful for large datasets
- YMMV

# pg\_partman

- Partition management in PostgreSQL
- Create partitions into the future before data arrives
- Archive/DROP older partitions for improved maintenance
- Use in conjunction with pg\_cron for automated creation and maintenance over time

# pg\_vector

- New extension adding vector features to PostgreSQL
- Useful to working with LLM embeddings
- Lots of focus and development happening currently across multiple vector extensions

# pg\_hint\_plan

- PostgreSQL doesn't support query hints or query plan caching
- Community extension that provides some ability to hint queries (similar to Oracle syntax)
- Won't solve all problems, but it can be helpful in a pinch

# pgcrypto

- TDE is not currently supported in vanilla PostgreSQL
- **pgcrypto** provides cryptographic functions to:
  - Apply symmetric/public key encryption of data
  - Work with OpenPGP keys and payloads
  - Generate passwords
  - Hash/checksum content

# TimescaleDB or Citus (BONUS)

- Examples of creating a "new" database on top of PostgreSQL using only extensions
- Automation of partitioning and sharding
- Various features for managing and querying large amounts of data across multiple nodes

# Monitoring and Instrumentation

# Server Metrics



- Error logs, off by default
- Cumulative Statistics System
  - Reset on crash, failover, or point in time restore
  - Can be manually reset
  - Aggregations only

# Cumulative Statistics System



- pg\_stat\_activity – current activity on server
- pg\_stat\_database – activity & size of database
- pg\_stat\_\*\_tables – activity & size of table
  - all
  - sys
  - user
- pg\_stat\_\*\_indexes – guess
- pg\_statio\_\*\_tables – specific info on blocks & storage

# Query Metrics



- EXPLAIN ANALYZE
- pg\_stat\_statements
  - Disabled by default

**DEMO!**



PASS Session

# Professional PostgreSQL

## Monitoring Made Easy

-  Friday, Nov. 8
-  2:30-3:30
-  Room 447



# Performance Tuning

# Cost-based Optimizer

- Statistics determine plan choice
- Customizable per server, table, and column
- Asynchronous process maintains statistics
- Manually update with ANALYZE

# Statistics

- Histogram of column values
- Defaults to 100 buckets
- Helpful views:
  - pg\_catalog.pg\_stats
  - pg\_catalog.pg\_stat\_user\_tables

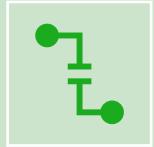
# EXPLAIN in Practice

- EXPLAIN = Estimated plan
- EXPLAIN ANALYZE = Actual plan
- EXPLAIN (ANALYZE,BUFFERS)
  - Query plan with disk IO
- EXPLAIN (ANALYZE,BUFFERS,VERBOSE)
  - Additional details on columns, schemas, etc.

# EXPLAIN

- Inverted tree representation
- There is no built-in visual execution plan
- EXPLAIN provides textual plan
  - pgAdmin and some others do attempt some visualizations
- Websites for visualizations and suggestions
  - <https://www.pgmustard.com/>
  - <https://explain.depesz.com/>

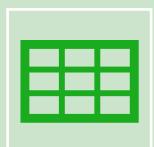
# Reading EXPLAIN



--> Nodes



Read inside-out



Join/Aggregate/Merge/Append at each level

```

EXPLAIN (ANALYZE, BUFFERS)
SELECT customer_name, count(*), date_part('year',order_date)::int order_year
FROM sales.orders o
    INNER JOIN sales.customers c ON o.customer_id=c.customer_id
WHERE c.customer_id=100
GROUP BY customer_name,order_year
ORDER BY order_year DESC;

```

GroupAggregate (cost=1827.91..1830.76 rows=102 width=35) (actual time=11.956..13.754 rows=4 loops=1)
 Group Key: ((date\_part('year'::text, (o.order\_date)::timestamp without time zone))::integer), c.customer\_name
 Buffers: shared hit=903
 -> Sort (cost=1827.91..1828.18 rows=107 width=27) (actual time=11.700..12.686 rows=107 loops=1)
 Sort Key: ((date\_part('year'::text, (o.order\_date)::timestamp without time zone))::integer) DESC, c.customer\_name
 Sort Method: quicksort Memory: 33kB
 Buffers: shared hit=903
 -> Nested Loop (cost=0.28..1824.30 rows=107 width=27) (actual time=0.113..10.649 rows=107 loops=1)
 Buffers: shared hit=903
 -> Index Scan using pk\_sales\_customers on customers c (cost=0.28..2.49 rows=1 width=27)
 Index Cond: (customer\_id = 100)
 Buffers: shared hit=3
 -> Seq Scan on orders o (cost=0.00..1819.94 rows=107 width=8) (actual time=0.053..8.413 rows=107 loops=1)
 Filter: (customer\_id = 100)
 Rows Removed by Filter: 73488
 Buffers: shared hit=900

Planning Time: 0.267 ms  
 Execution Time: 13.934 ms

GroupAggregate (**cost=1827.91..1830.76 rows=102 width=35**)  
**(actual time=11.956..13.754 rows=4 loops=1)**  
Group Key: ((date\_part('year'::text, (o.order\_date)::timestamp without time zone))::integer), c.customer\_name  
**Buffers: shared hit=903**  
-> Sort (**cost=1827.91..1828.18 rows=107 width=27**) (**actual time=11.700..12.686 rows=107 loops=1**)  
    Sort Key: ((date\_part('year'::text, (o.order\_date)::timestamp without time zone))::integer) DESC, c.customer\_name  
**Sort Method: quicksort Memory: 33kB**  
Buffers: shared hit=903  
-> **Nested Loop** (**cost=0.28..1824.30 rows=107 width=27**) (**actual time=0.113..10.649 rows=107 loops=1**)  
    Buffers: shared hit=903  
        -> **Index Scan** using pk\_sales\_customers on customers c (**cost=0.28..2.49 rows=1 width=27**)  
            Index Cond: (customer\_id = 100)  
            Buffers: shared hit=3  
        -> **Seq Scan** on orders o (**cost=0.00..1819.94 rows=107 width=8**) (**actual time=0.053..8.413 rows=107 loops=1**)  
            Filter: (customer\_id = 100)  
            Rows Removed by Filter: 73488  
            Buffers: shared hit=900  
Planning Time: 0.267 ms  
Execution Time: 13.934 ms

# Primary Scan nodes

- Sequence scan (seq scan)
- Index scan
- Index-only scan

# Primary join nodes

- Nested Loop
- Hash Join
- Merge Join

Explain Glossary: <https://www.pgmustard.com/docs/explain>

# Other nodes

- Aggregate
- Append
- Sort
- Group
- Window Aggregate
- More...

Explain Glossary: <https://www.pgmustard.com/docs/explain>



# `pg_stat_statement`

- MUST HAVE Extension for query analysis
- Installed in all/most DBaaS offerings
- Does not save plans or provide recommendations
- Allows you to do triage: “Is this a big query issue, or death by 1,000 small cuts?”

# **pg\_stat\_statement statistics**

- Execution Time (total/min/max/mean/stddev)
- Planning Time (total/min/max/mean/stddev)
- Calls (total)
- Rows (total)
- Buffers (shared/local/temp)
  - read hit dirtied written
  - read/write time
- WAL

43 Columns of data  
(as of PG16)

Name	Value
userid	16422
dbid	16434
queryid	-6155333619461995114
query	SELECT id, name FROM...
plans	0
total_plan_time	0.0
min_plan_time	0.0
max_plan_time	0.0
mean_plan_time	0.0
stddev_plan_time	0.0
calls	151
total_exec_time	8.489053
min_exec_time	0.013751
max_exec_time	1.356096
mean_exec_time	0.056218894039735096
stddev_exec_time	0.11851139585068957
rows	151
shared_blk_hit	450
shared_blk_read	3

All statistics are cumulative  
from the last restart\*

\*or reset by a superuser



# pg\_stat\_statement: total time w/ cache hit ratio

```
SELECT query, calls, total_time, rows,  
       100.0 * shared_blk_hit/nullif(shared_blk_hit + shared_blk_read, 0)  
          AS hit_percent  
FROM pg_stat_statements  
ORDER BY total_time DESC  
LIMIT 5;
```

```
-[ RECORD 1 ]-----  
query           | SELECT * FROM sales.orders where customer_id = ?  
calls           | 1000  
total_time      | 2451.000000000000  
Rows.            | 521  
hit_percent     | 99.9873421
```

\*Lots of good videos/tutorials, but also [see the docs](#) for more help and ideas.

PASS Session

# Explaining the Postgres Query Planner



Wednesday, Nov. 6



2:00-3:00



Room 445-446



PASS Session  
**Can SQL Server  
DBAs/Developers Tune  
Queries in PostgreSQL**

-  Friday, Nov. 8
-  10:30-11:30
-  Room 447



# PostgreSQL and the Cloud

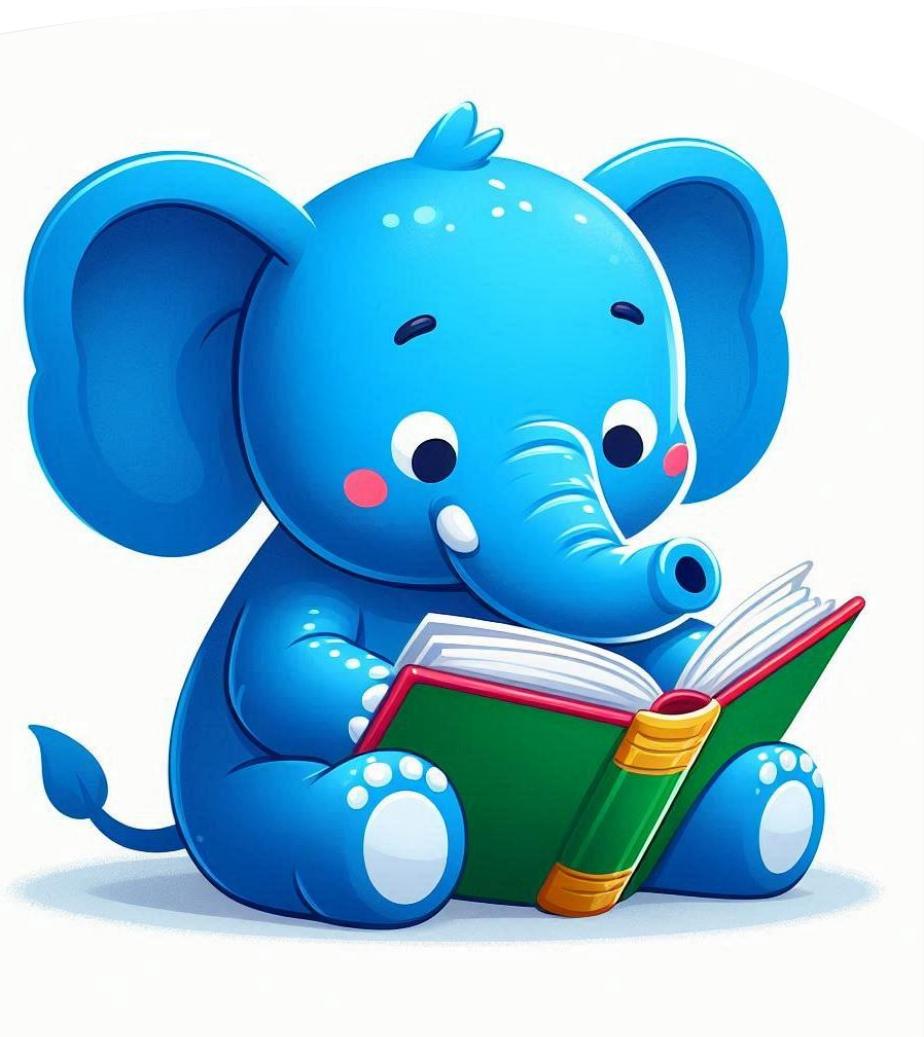
PASS Session  
**PostgreSQL in the Cloud**

 Thursday, Nov. 7  
 9:45-10:45  
 Room 445-446



# Additional Resources

# PostgreSQL Documentation



- [https://www.postgresql.org/docs/](https://www.postgresql.org/docs/current/dml-delete.html)  
cs/
- Pay attention to this:
  - [https://www.postgresql.org/docs/c  
urrent/dml-delete.html](https://www.postgresql.org/docs/current/dml-delete.html)

# Books



- OURS!
- [Art of PostgreSQL](#) – Dimitri Fontaine
- [Database Administration](#) – Craig Mullins
- [PostgreSQL Query Optimization](#) - Henrietta Dombrovskaya, Boris Novikov and Anna Baillieкова

# Other Events



- [PostgreSQL User Groups \(PUG\)](#)
- [MeetUp](#)
- [PGDay](#)
- pgConf (same link as PGDay)

# Online



- [Planet PostgreSQL](#)
- [PostgreSQL Slack](#)
- Cooper Press [Email List](#)
- #pghelp on X
- Usual suspects

# Calling All PostgreSQL Users

**Whether you're a pro or just getting started,  
Redgate would love to hear your PostgreSQL  
insights and pain points**

Together, we'll co-design the future of  
PostgreSQL tools and make development  
smoother for everyone!

