

Pgloader Overview and Features

what is Pgloader ? how it works internally ? what are the characteristics of pgloader ? what are the features it has ?

explain all about pgloader so that i donot have any doubt in this and have cristal and clear knowledge of pgloader

Below is a **complete, crystal-clear, zero-doubt explanation of pgloader**, written in **simple but deep technical language**, step-by-step.

I'll explain **what it is** → **why it exists** → **how it works internally** → **characteristics** → **features** → **real-world behavior** → **when to use & when not**.

1 What is pgloader?

pgloader is an open-source data migration tool used to **load data into PostgreSQL** from other databases, files, or formats.

👉 In simple words:

pgloader helps you move data from one system into PostgreSQL safely, fast, and automatically.

Common use cases

- MySQL → PostgreSQL
- MSSQL → PostgreSQL
- SQLite → PostgreSQL
- CSV / TSV files → PostgreSQL
- Old legacy DB → PostgreSQL

2 Why pgloader exists (The core problem it solves)

Migrating databases is **hard**, because:

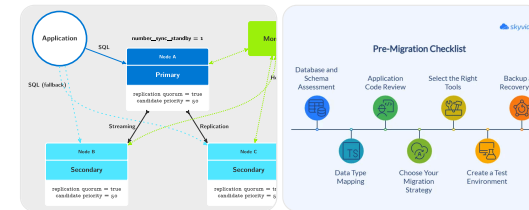
- Different databases have **different data types**
- Constraints behave differently

- Indexes and foreign keys cause failures
- Large datasets are **slow** to migrate
- Manual scripts break easily

pgloader solves this by:

- ✓ Automatically converting data types
- ✓ Loading data in **parallel** (very fast)
- ✓ Handling constraints smartly
- ✓ Retrying, logging, and validating

3 High-level working flow of pgloader



Simple flow

```
pgsql

Source Database/File
↓
pgloader Engine
↓
PostgreSQL Database
```

4 How pgloader works internally (DEEP EXPLANATION)

This is the most important part 📌

I'll explain it like **inside-the-engine behavior**.

♦ Step 1: Source Connection

pgloader first connects to the **source system**:

- MySQL / MSSQL / SQLite
- Or file system (CSV, TSV)

It reads:

- Tables
- Columns
- Data types
- Constraints
- Index definitions

👉 This metadata is stored **internally in memory**.

♦ **Step 2: Schema Analysis & Type Mapping**

pgloader **does not blindly copy data**.

Instead, it:

- Maps source data types → PostgreSQL types

Example:

Source DB	Source Type	PostgreSQL Type
MySQL	INT	INTEGER
MySQL	TINYINT(1)	BOOLEAN
MySQL	DATETIME	TIMESTAMP
MSSQL	NVARCHAR	TEXT

This mapping is:

- ✓ Automatic
- ✓ Configurable (you can override it)

♦ **Step 3: Target Preparation (VERY IMPORTANT)**

Before inserting data, pgloader prepares PostgreSQL:

- ✓ Creates tables
- ✓ Creates sequences
- ✓ Temporarily **disables constraints**

- Foreign keys
- Triggers
- Indexes

👉 Why?

Because inserting millions of rows with constraints ON is **slow and error-prone**.

♦ **Step 4: Parallel Data Loading (Key Performance Feature)**

pgloader loads data using:

- **Multiple worker threads**
- **COPY command** (fastest PostgreSQL insertion)

Internally:

```
mathematica

Thread 1 → Table A
Thread 2 → Table B
Thread 3 → Table C
```

👉 This is why pgloader is **much faster than manual INSERT scripts**.

♦ **Step 5: Data Transformation (On-the-fly)**

During loading, pgloader can:

- Convert encodings
- Fix invalid values
- Apply expressions
- Rename columns
- Skip bad rows (optional)

Example:

```
lisp
```

```
CAST type datetime TO timestamptz
```

◆ Step 6: Post-Load Cleanup

After data is loaded:

- ✓ Re-enable constraints
- ✓ Create indexes
- ✓ Validate foreign keys
- ✓ Analyze tables (`ANALYZE`)

Now PostgreSQL is **ready for production use**.

5 Characteristics of pgloader (Core Nature)

These define **how pgloader behaves**.

◆ 1. PostgreSQL-Focused

- Built **specifically** for PostgreSQL
- Uses PostgreSQL internals like:
 - COPY
 - ANALYZE
 - Sequences

◆ 2. Declarative Configuration

You describe **what you want**, not how.

Example:

```
sql

LOAD DATABASE
  FROM mysql://user:pass@localhost/db
  INTO postgresql://user:pass@localhost/db;
```

◆ 3. Fault-Tolerant

- Can skip errors
- Logs rejected rows
- Supports retries

◆ 4. High Performance

- Parallel loading
- Bulk COPY
- Constraint optimization

◆ 5. Non-Destructive

- Does NOT modify source DB
- Only reads from it

6 Features of pgloader (What it can do)

✓ Database to PostgreSQL

- MySQL → PostgreSQL
- MSSQL → PostgreSQL
- SQLite → PostgreSQL

✓ File-based Loading

- CSV
- TSV
- Fixed-width files

✓ Automatic Schema Creation

- Tables
- Sequences
- Indexes

✓ Data Type Conversion

- Automatic mapping
- Custom rules

✓ Constraint Handling

- Disable → load → enable
- Avoids FK failures

✓ Parallelism

- Multiple tables at once
- Massive speed boost

✓ Encoding Handling

- UTF-8 fixes
- Character set conversion

✓ Logging & Monitoring

- Progress output
- Error reports
- Rejected rows file

✓ Scriptable & Repeatable

- Same script can be run again
- Ideal for CI/CD migrations

7 Real-World Example (Conceptual)

MySQL → PostgreSQL migration

```
pgsql
1 million users
50 tables
Foreign keys everywhere
```

Without pgloader

- ✗ Manual scripts
- ✗ FK errors
- ✗ Takes hours or days

With pgloader

- ✓ 1 command
- ✓ Handles FK automatically
- ✓ Finishes in minutes

8 When should you use pgloader?

✓ Use pgloader when:

- ✓ Migrating to PostgreSQL
- ✓ Large datasets
- ✓ Schema + data both needed
- ✓ Speed matters

✗ Do NOT use pgloader when:

- ✗ PostgreSQL is NOT the target
- ✗ You need real-time replication
- ✗ You need bi-directional sync

(pgloader is for **migration**, not replication)

9 pgloader vs Liquibase (Important clarity)

pgloader	Liquibase
Data migration	Schema versioning
One-time bulk load	Continuous DB evolution
External sources	PostgreSQL only
High-volume data	Small incremental changes

👉 They solve different problems.

10 Mental Model (Remember this)

pgloader is a migration engine that reads, transforms, and bulk-loads data into PostgreSQL with speed and safety.

If you understand:

- ✔ Source reading
- ✔ Type mapping
- ✔ Parallel COPY
- ✔ Constraint management

👉 You fully understand pgloader.

If you want, next I can:

- Explain pgloader vs AWS DMS
- Show real pgloader scripts
- Explain error handling scenarios
- Explain pgloader in production migration

Just tell me 👍

how thread works in pgloader

Below is a **deep, crystal-clear explanation of how threads work in pgloader**, written so that **you won't have any confusion**, even internally/architecturally.

I'll go step-by-step and also give you a **mental model** you can always remember.

1 What does “thread” mean in pgloader?

In **pgloader**, a **thread** is an **independent worker** that performs part of the migration in **parallel**.

👉 Each thread:

- Reads data from the source
- Transforms it (if needed)
- Loads it into PostgreSQL using `COPY`

So instead of:

```
CSS
Table A → done
Table B → done
Table C → done
```

pgloader does:

```
mathematica
Table A → Thread 1
Table B → Thread 2
Table C → Thread 3
```

This is the **core reason pgloader is fast**.

2 Why pgloader needs threads (the problem it solves)

Without threads:

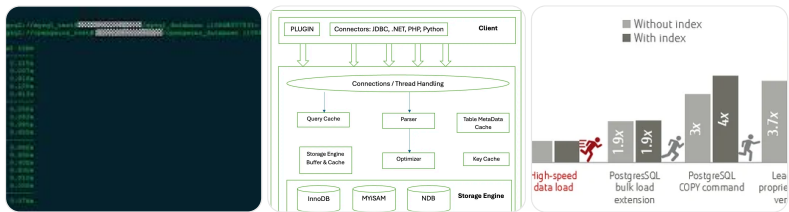
- One table loads at a time
- CPU is mostly idle
- Network is underused

- Migration is slow

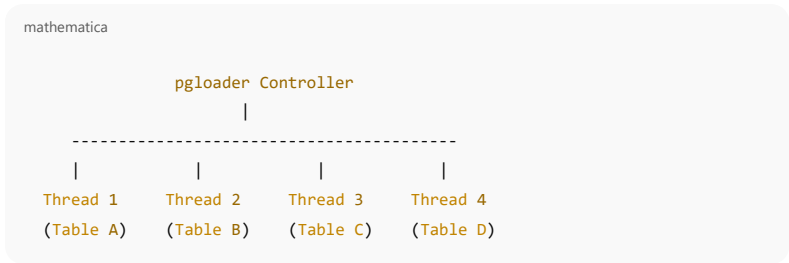
With threads:

- ✓ CPU fully utilized
- ✓ Network fully utilized
- ✓ Disk IO parallelized
- ✓ PostgreSQL COPY speed maximized

3 High-level threading architecture



Think of pgloader like this:



Each thread is independent but coordinated.

4 Internal thread lifecycle (VERY IMPORTANT)

Let's follow one pgloader thread from birth to death.

♦ Step 1: Thread creation

pgloader:

- Detects how many tables need to be loaded
- Creates a **thread pool**

Default:

```
ini
workers = number of CPU cores
```

(You can override this.)

♦ Step 2: Task assignment

Each thread is assigned a **task**, such as:

- Load one table
- Load part of a large table (chunking)

Example:

```
mathematica
Thread 1 → users
Thread 2 → orders
Thread 3 → products
```

♦ Step 3: Source reading (per thread)

Each thread:

- Opens its **own source connection**
 - Reads rows in batches
 - Streams data (does NOT load entire table in memory)
- 👉 This avoids memory explosion.

♦ Step 4: Data transformation (inside thread)

Inside each thread:

- Data types are converted

- NULLs fixed
- Encoding corrected
- Invalid values handled

This happens **row-by-row, inside the thread**.

◆ Step 5: COPY into PostgreSQL

This is the most critical part.

Each thread:

- Opens its **own PostgreSQL connection**
- Executes:

```
sql
COPY table_name FROM STDIN;
```

So internally:

```
powershell

Thread 1 → COPY users
Thread 2 → COPY orders
Thread 3 → COPY products
```

👉 PostgreSQL handles multiple COPY operations **concurrently**.

◆ Step 6: Thread completion

After finishing:

- Thread reports success/failure
- Releases DB connections
- Returns to thread pool (or exits)

5 Two levels of parallelism in pgloader

This is where most people get confused.

◆ Level 1: Table-level parallelism

Different tables loaded at the same time.

```
mathematica

users → Thread 1
orders → Thread 2
items → Thread 3
```

◆ Level 2: Chunk-level parallelism (large tables)

If a table is **huge**, pgloader can:

- Split it into chunks
- Load chunks in parallel

Example:

```
mathematica

users (10M rows)
├─ chunk 1 → Thread 1
├─ chunk 2 → Thread 2
└─ chunk 3 → Thread 3
```

6 How pgloader avoids thread conflicts

Parallelism is dangerous if unmanaged. pgloader avoids problems by:

✓ Constraint handling

- Foreign keys disabled before loading
- Re-enabled after all threads finish

✓ Index handling

- Indexes created **after** data load
- No index contention during COPY

✓ Transaction isolation

- Each thread uses its own transaction

- No shared locks on source

7 Thread configuration in pgloader

You can control threads explicitly.

Example:

```
lisp

WITH
  workers = 4,
  concurrency = 2
```

Meaning:

- `workers` → total worker threads
- `concurrency` → how many tables at the same time

8 What happens if threads are too many?

✗ Too many threads causes:

- PostgreSQL connection exhaustion
- Disk IO saturation
- Slower migration (counter-productive)

✓ Best practice:

```
nginx

threads = CPU cores
```

9 Failure handling with threads

If one thread fails:

- Other threads continue
- Error is logged
- Migration does NOT crash completely (unless fatal)

This is called **partial fault tolerance**.

10 Mental model (remember forever)

pgloader threads are independent workers that read, transform, and COPY data in parallel, while a central controller coordinates constraints and order.

If you understand:

- ✓ Thread pool
- ✓ One thread = one COPY stream
- ✓ Parallel tables / chunks
- ✓ Constraints disabled globally

👉 You have **complete mastery of pgloader threading**.

11 Short real-world analogy

Imagine:

- 1 warehouse (PostgreSQL)
- 5 trucks (threads)
- Each truck brings goods (data)
- Security gates (constraints) are opened **after all trucks arrive**

That's exactly pgloader threading.

If you want next, I can:

- Explain **workers vs concurrency** with diagrams
- Explain **pgloader performance tuning**
- Explain **thread behavior during FK-heavy schemas**
- Compare **pgloader threads vs Java threads**

Just tell me 🙌