

Assignment 1

Uploading Data to the Database

Database Tuning

Start date: October 16, 2014

Due date: October 28, 23:59

Grading: 1 point

Note: This assignment involves reading the documentation of Java, the JDBC driver, and PostgreSQL. Finding the relevant sources of information is part of the challenge.

1. Download `dblp.zip` at:

<http://www.cosy.sbg.ac.at/~augsten/teaching/2014ws/dbt/dblp.zip>

This archive contains two tab separated files (`publ.tsv` and `auth.tsv`) that store authors and their publications as found in the DBLP¹ bibliography. The imported tables have the following schemas:

- `Auth(name(49),pubID(129))`
- `Publ(pubID(129),type(13),title(700),booktitle(132),year(4),publisher(196))`

You can assume that all attribute values are strings; the maximum string length is shown in brackets.

2. The straightforward algorithm to load the data from the TSV file to a table issues an SQL INSERT query for each line in the TSV file.

Task 1: Implement the straightforward approach to load `auth.tsv` to the database (PostgreSQL, Java).²

Task 2: The straightforward approach is slow. There are other approaches that are significantly faster. Figure out how the efficient approaches work and implement two of them.

Report: Describe the two efficient approaches. Give the runtime for loading `auth.tsv` with the straightforward and the efficient approaches. Why are the efficient approaches faster? Which tuning principle did you apply?

Access parameters for PostgreSQL:

host: `dumbo.cosy.sbg.ac.at`, port: 5432

¹<http://www.informatik.uni-trier.de/~ley/db>

²This might be very slow. Instead of loading all the data with this approach, you can also load part of the data and assume that runtime scales linearly. Mention this in the report.

db: `pstuningws1415`

user/password: you should have received them via email

The database server (`dumbo`) is accessible only from inside the university network. If you would like to work from home, please connect to `fanny.cosy.sbg.ac.at` via ssh. Java and the PostgreSQL client are installed on this machine.

Please indicate the average time per group member that was spent solving this assignment. The time that you indicate will have *no* impact on your grade.