



Assignment 3

Hand in this assignment until **Thursday, 18 May 2023, 12:00** at the latest.

Exam-style Exercises

Exercises marked with **E** are similar in style to those you will find in the exam. You can use these to hone your expectations and gauge your skills.

Running out of ideas?

Are you hitting a roadblock? Are some of the exercises unclear? Do you just need that one hint to get the ball rolling? Refer to the [#forum](#) channel on our Discord server and check the tag for this assignment—maybe you'll find just the help you need.

Task 1: Declarativity and Query Execution Strategies **E**

(1 credit)

Database queries are written in a *declarative manner*: The user specifies “*what*” the desired result is, but not “*how*” to compute it. But before we delve into this declarative world of SQL and let the query processor of PostgreSQL do the job, we are going to explore the effects that different – but semantically equivalent – execution strategies can have on query runtime.

We will query the LEGO dataset (`contains.csv`, `bricks.csv`, `minifigs.csv` and `sets.csv`) which should be familiar from the lecture.

The following naive PyQL query returns all minifigs of category '67' that are part of sets '671-1' or '377-1'. However the query does not perform well:

```
1 from DB1 import Table
2
3 contains = Table('contains.csv')
4 bricks   = Table('bricks.csv')
5 minifigs = Table('minifigs.csv')
6
7 for c in contains:
8     for m in minifigs:
9         if c['piece'] == m['piece']:
10            if (c['set'] == '671-1' or c['set'] == '377-1'):
11                if m['cat'] == '67':
12                    print({ 'set': c['set'], 'name': m['name'] })
```

Implement an alternative version that is *as efficient as possible* by exploiting properties of the data and the query itself, such as:

Order of rows, early filtering of data, uniqueness or disjointness of values in columns, query strategy, etc. (Keep in mind that not all of them may apply here.)

Note

As in the lecture, you may use the loop control keywords `break` (and `continue`) or `for ... else`.

Task 2: Query Optimization through Data Pre-Processing

(1 credit)

In the lecture, we discussed multiple variations of pre-processing the data to further improve the efficiency of the “weight of LEGO set 5610” query. Another such strategy was proposed by a fellow student in a previous course.

The idea is simple: A **piece** referenced in **contains.csv** may be found in either **bricks.csv** or **minifigs.csv**. To look up a piece’s details, we have to search both files since we can’t tell from a piece identifier whether it refers to a brick or a minifig. However, this knowledge could be explicitly represented in the data by extending **contains.csv** with an additional column **type**: For each entry in **contains.csv**, **type** indicates whether the piece is a brick ('b') or minifig ('m').

- (a) Construct a new file **contains_type.csv**¹, which is a copy of **contains.csv** but with an additional column **type** as described above. Approach this problem in two steps:

- i. Write a PyQL query that constructs the rows to be placed in **contains_type.csv**. In this first step, the query will generate a Python list **l** of dictionaries.

Hint

A naive query would probably take far too long (~ several hours). Instead, try to optimize your PyQL query. Consider that **minifigs.csv** is a small table compared to **bricks.csv** (~ 17% of all pieces). We also know that if a piece is not a minifigure, it is a brick. Further, pieces are unique and **minifigs.csv** is sorted by its column **piece**.

Note

Even when optimized your query may take some minutes to process. So think carefully about possible improvements and be patient while running your query.

- ii. Use method **dump** from **DB1.py** to construct a Table object from the list **l** and dump it to a file in the usual CSV format.

Consider the following example:

```
1 l = [{ 'a': 23, 'b': 'foo' }, { 'a': 42, 'b': 'bar' }]
2 Table(l).dump('foo.csv')
```

- (b) Reimplement the “weight of LEGO set 5610” query leveraging the new **type** column. Compare the query execution times for the lecture’s baseline query (slide 12 in slide set **db1-03.pdf**) and your new optimized variant.
- (c) Provide a formula that estimates the work performed by the optimized PyQL query. (For reference, see slide 13 and 15 in slide set **db1-03.pdf** about *Declarativity*.)

¹This file has also been provided to you in the **legosets.zip** archive.