

Assignment

A Query Engine for Acyclic Conjunctive Queries

Data and Information Management



Professor Bas Ketsman <bas.ketsman@vub.be>
Assistant Tim Baccaert <tim.baccaert@vub.be>

The deadline for this assignment is **Friday 29th, December 2023** at 23:59, Central European Time (GMT+1). You have to upload your solution to Canvas, as prescribed in the third section of this document. This assignment is designed to be made in **pairs**, if this poses a problem you may work **individually** as well.

I. Context

In this assignment, your goal is to implement a query engine for **acyclic conjunctive queries**, using the techniques and concepts we encounter during the lectures. This assignment is mandatory and constitutes your grade for the practical segment of this course, which is 40% of your overall grade.

We designed the assignment for **pairs** of students. You are free to assemble your group as you see fit, you do not need to notify us about its composition. In case you have trouble finding suitable colleagues for whatever reason, you may alternatively work on this assignment **individually**. We will, however, not grant you extra points if you do choose to work alone.

Programming Language

Since this is a course at the master's level, we have to accomodate students from a wide range of different universities. However, we assume you have advanced knowledge of at least one programming language, as can be reasonably expected from someone with a bachelor's degree in Computer Science. To facilitate this difference in background, we allow you to choose from the following selection of languages for this assignment:

- Java (<https://dev.java>)
- Rust (<https://www.rust-lang.org>)
- Scala (<https://scala-lang.org>)
- TypeScript (<https://www.typescriptlang.org>)

If you are not familiar with any of these languages, you may request an exemption by sending an email to tim.baccaert@vub.be. We strongly prefer the use of a statically typed language. Furthermore, the language must have a library with bindings to Apache Arrow (<https://arrow.apache.org>), as this assignment depends on it.

You are allowed to use libraries, provided they do not implement required algorithms for your assignment (i.e., you may use libraries for logging, parsing, data structures, etc.; but not join algorithms, query engines, ...).

Note: We will not answer questions regarding the specifics of a programming language, or used libraries. You should be mature enough as a programmer to figure out these concerns yourself. Questions should exclusively relate to the concepts required for the assignment, or assignment practicalities.

Testing Dataset

Throughout this assignment we will make use of the dataset accessible from the Canvas space for testing purposes. It is comprised of five files in the comma-separated values (CSV) format, each file represents a single relation. A visual representation of the schema is given in Figure 1.

The relation **Beers** (`beers.csv`) has 8 attributes: `beer_id`, the artificial primary key of the beer; `brew_id`, the id of its brewery; `beer`, its name; `abv`, its alcohol by volume; `ibu`, its international bitterness unit; `ounces`, its typical volume measured in ounces; and finally `style` and `style2`, its brewing styles.

Beers has a foreign key relationship with the **Styles** (`styles.csv`) relation, which has 3 attributes: `style_id`, an artificial primary key (shared with `cat_id`) for the style; `cat_id`, the id of its category; and `style`, the textual label for this style. This relation is in a foreign key relationship with **Categories** (`categories.csv`), which has 2 attributes `cat_id` and `cat_name`, respectively its artificial primary key and category name.

Next, we have the **Locations** (`locations.csv`) relation, with the 3 attributes: `loc_id` and `brew_id` forming a shared primary key, where the latter is the id of the brewery; `longitude` and `latitude` representing a brewery's geographic location; and `accuracy` representing the accuracy of its location.

Lastly, we have the **Breweries** (`breweries.csv`) relation, which has the following 11 attributes: `brew_id`, the artificial primary key of the brewery; `brew_name`, the name of the brewery; and the self-explanatory attributes `address1`, `address2`, `city`, `state`, `code`, `country`, `phone`, `website`, and finally `description`.

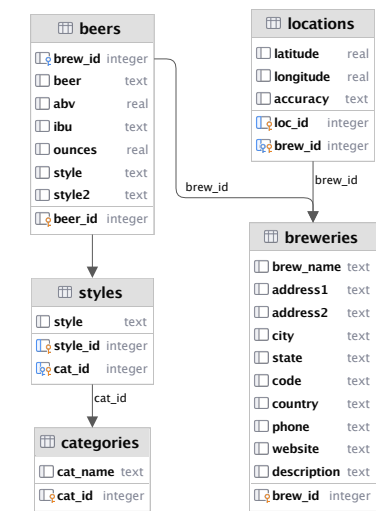


Fig 1. visual schema representation.

Apache Arrow

To store relations in memory, we expect you to make use of Apache Arrow (<https://arrow.apache.org>). This library defines a language-independent format for in-memory storage of relational data, optimised for analytic queries on modern hardware. Think of it as a portable in-memory database without a query engine and transaction support. Arrow supports loading data from many file formats, and we will use CSV files for this assignment, as mentioned previously.

The concrete API varies from language to language, hence we refer you to the corresponding websites for each of the languages: Java (<https://arrow.apache.org/docs/java/index.html>), Rust (<https://docs.rs/arrow/latest/arrow/>), Scala (see Java), and Typescript (<https://www.npmjs.com/package/@apache-arrow/ts>). Installation depends on your language environment, we suggest you use your language's standard package manager and repositories (*i.e.*, `sbt` for Scala, `cargo` for Rust, `npm` for Typescript).

II. Assignment

In the following subsections we briefly go over the features you are expected to implement. Each feature has a number of fixed requirements, and may include bonus requirements that you can use to distinguish yourself.

We only give points for bonus requirements if you implemented all the regular requirements.

Part I. Query Representation

In this assignment, we focus on conjunctive queries (CQs). This is a practically relevant fragment as most user queries are CQs. As seen in the lectures, we use datalog notation.

F1. Data Structure

First, we expect you to implement a data structure that syntactically represents a CQ. It has the following components and subcomponents:

- a **head** atom *e.g.*, $\text{Answer}(\bar{y})$,
- a set of **body** atoms *e.g.*, $R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$,
- each **atom** consists of a *relation name*, and a tuple of *terms* *e.g.*, $R_1(u_1^1, \dots, u_1^k)$, and
- each **term** u_j^i is either a *constant value* or a *variable*.

Here, a constant value is a value that has one of the Apache Arrow logical types¹: Utf8, FloatingPoint, and Int. Depending on the language you are using, these may have a slightly different name; use the most similar ones. The precise data structures you use to represent these components are up to you, and will likely vary based on the rest of your implementation.

Example 1. The following query asks if there are beers that are categorised as Belgian and French ales:

```
Answer() :- Beers(beer_id, brew_id, beer, abv, ibu, ounces, style, style2),
           Styles(style_id, cat_id, style),
           Categories(cat_id, 'Belgian and French Ale').
```

The head atom is `Answer()` and the body atoms are `Beers(...)`, `Styles(...)`, and `Categories(...)`. All lowercase symbols such as `beer` and `style` are variable terms, and `'Belgian and French'` is a constant value of type `Utf8`.

Requirement 1 Implement a Query data structure.

As bonus functionality, you can implement a parser for the rule syntax.

(Bonus Functionality) Implement parsing functionality according to the Backus-Naur form (BNF) grammar below. It should produce an instance of your Query data structure (you may use a parsing library).

```
<query> ::= <atom> ":-" <atom-list> "."
<atom> ::= <relation> "(" <term-list> ")"
<atom-list> ::= <atom> | <atom> "," <atom-list> | ""
<term> ::= <constant> | <variable>
<term-list> ::= <term> | <term> "," <term-list> | ""
<constant> ::= <float> | <integer> | <string>
<variable> ::= <letters>
<relation> ::= <capital> <letters>
<float> ::= <digit> "." <digits> | "-" <digit> "." <digits>
<integer> ::= <digit> | "-" <digit> | <digits>
<string> ::= "'" <letters> "'"
<capital> ::= "A" | .. | "Z"
<letter> ::= <capital> | "a" | .. | "z"
<digit> ::= "0" | .. | "9"
<letters> ::= <letter> <letters> | ""
<digits> ::= <digit> <digits> | ""
```

Lst 1. conjunctive query BNF grammar

¹<https://github.com/apache/arrow/blob/4fc2731a3ceaf7dd5b1dce6f29bf7cad0ab2f13e/format/Schema.fbs#L407-L430>

F2. Testing for Acyclicity

Since we will specifically focus on **acyclic** conjunctive queries, we need to make sure that the query we are given is *acyclic*. Recall from the lectures that a CQ is acyclic if the Graham-Yu-Ozsoyoglu (GYO) algorithm yields the empty query.

Example 2. Given the query from Example 1, the GYO algorithm should produce the empty query. A possible execution removes the following ears:

1. the ear Beers(...) with witness Styles(...),
2. the ear Styles(...) with witness Categories(...), and
3. the ear Categories(...), which has no shared variables, and is therefore its own witness.

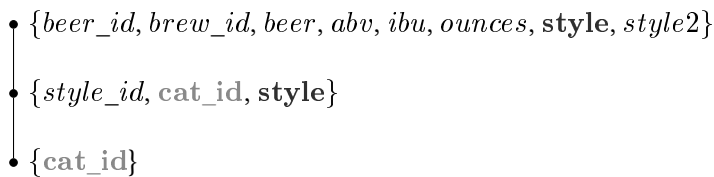
After removing these ears, the query is empty, hence it is acyclic.

Requirement 2 Implement the GYO algorithm for your Query data structure.

F3. Join Trees

From the lectures, we know that a CQ is acyclic only if it has a join tree. Hence, for the next feature we ask you to implement a procedure that can construct such a join tree based on GYO's execution.

Example 3. A possible join tree based on the GYO execution in Example 2 is the following:



Requirement 3 Implement an algorithm that constructs a join tree for your Query data structure.

Part II. Query Evaluation

We now move on to the evaluation of acyclic CQs. Here, we rely on the Yannakakis algorithm which was handled in the lectures. Note that in order for this algorithm to work, the query **must** be acyclic.

F4. Globally Consistent Databases

A globally consistent database is a database that does not contain any dangling tuples. We can use the join tree to build a **full reducer** for our database. When applied over the database, this yields us a globally consistent version.

Example 4. Given that the node with atom Beers(..) is chosen as the root. A possible application of a full reducer based on the join tree from Example 3 is given by:

$$\begin{aligned} \text{Styles}' &= \text{Styles} \times \text{Categories} && (\text{post-order}) \\ \text{Beers}' &= \text{Beers} \times \text{Styles}' && (\text{post-order}) \\ \text{Styles}'' &= \text{Styles}' \times \text{Beers}' && (\text{pre-order}) \\ \text{Categories}' &= \text{Categories} \times \text{Styles}'' && (\text{pre-order}) \end{aligned}$$

Requirement 4 Implement an algorithm that creates a globally consistent database given an input database and a join tree.

F5. The Yannakakis Algorithm

Finally, you will bring all components together into a single algorithm. First, verify whether the query given to you by the user is acyclic, and terminate with an error if it is cyclic. Then, create a join tree for the query, and apply the Yannakakis algorithm.

Example 5. Given the join tree from Example 3 and the globally consistent database obtained from applying the full reducer in Example 4, an evaluation of Yannakakis executes the following expression:

$$\pi_{()}(\text{Beers}' \bowtie (\text{Styles}' \bowtie \text{Categories}'))$$

The answer based on the test database should yield $\{\}$ (i.e., false) as the output.

Requirement 5 Implement the (generalized) Yannakakis algorithm for acyclic CQs.

(F6.) Additional Bonus Features

If you wish to distinguish yourself, you can extend your implementation with other algorithms regarding conjunctive queries. This could be a worst-case optimal join algorithm for cyclic CQs, it could include CQ minimisation, constant-time enumeration of CQ answers with linear preprocessing, ranked enumeration of CQ answers, CQ evaluation under inclusion or functional dependencies, and so on. You can look at the book for further inspiration.

III. Delivery

The deadline for the delivery of this project is on **Monday 29th, December 2023** at 23:59, Central European Time (GMT+1). Your report and implementation should be submitted in their respective Canvas locations.

Report You are expected to write a report that accompanies your assignment. Make use of the following ACM L^AT_EX-template for your report: <https://www.acm.org/publications/proceedings-template>, and use the `acmart` document class with the options `nonacm`, `sigconf`, and `screen`. The report is limited to a maximum of 4 pages excluding code listings, figures, and references.

Make sure to mention your own name (and that of your colleague, if there is one), including your student number and email. **Both students in a pair should submit the same assignment to Canvas.**

The report should be well-structured, and written in English. We expect you to discuss your implementation strategy for each of the requirements, and a short motivation for why you made your design decisions the way you did. **Name the file <last-name-1>-<last-name-2>-report.pdf or <last-name>-report.pdf.**

Implementation We will automate some aspects of the grading process. Hence, you should run your implementation on the following five conjunctive queries:

1. $\text{Answer}() :- \text{Beers}(u_1, x, u_2, '0.07', u_3, u_4, y, u_5), \text{Styles}(u_6, z, y), \text{Categories}(z, u_7), \text{Locations}(u_8, x, u_9, u_{10}, u_{11}), \text{Breweries}(x, u_{12}, u_{13}, u_{14}, u_{15}, u_{16}, u_{17}, u_{18}, u_{13}, u_{14}, u_{15}).$
2. $\text{Answer}(x, y, z) :- \text{Breweries}(w, x, 'Westmalle', u_1, u_2, u_3, u_4, u_5, u_6, u_7, u_8), \text{Locations}(u_9, w, y, z, u_{10}).$
3. $\text{Answer}(x, y, z) :- \text{Beers}(u_1, u_2, z, u_3, u_4, u_5, x, u_6), \text{Styles}(u_7, y, x), \text{Categories}(y, z).$
4. $\text{Answer}(x, y, z, w) :- \text{Beers}(u_1, v, x, '0.05', '18', u_2, 'Vienna Lager', u_3), \text{Locations}(u_4, v, y, z, w).$
5. $\text{Answer}(x, y, z, w) :- \text{Beers}(u_1, x, u_2, '0.06', u_3, u_4, y, u_5), \text{Styles}(u_6, z, y), \text{Categories}(z, w), \text{Locations}(u_8, x, u_9, u_{10}, u_{11}), \text{Breweries}(x, u_{12}, u_{13}, u_{14}, u_{15}, u_{16}, u_{17}, u_{18}, u_{13}, u_{14}, u_{15}).$

Your answers should be stored in a CSV file with the name `output.csv` which has the following columns:

1. **query_id** (Int): the identifier of the query (e.g., 1, 2, 3, ...).
2. **is_acyclic** (Int): whether the query is acyclic (1) or cyclic (0) after executing your implementation of the GYO algorithm.
3. **bool_answer** (Int) (when the query is boolean) whether the query outputs true (1) or false (0) after executing your implementation of the Yannakakis algorithm. Leave this NULL when the query is not boolean.
4. **attr_x_answer** all constants which get mapped to by x by all valuations for the query. (i.e., the output of the x column)
5. **attr_y_answer** all constants which get mapped to by y by all valuations for the query. (i.e., the output of the y column)
6. **attr_z_answer** all constants which get mapped to by z by all valuations for the query. (i.e., the output of the z column)
7. **attr_w_answer** all constants which get mapped to by w by all valuations for the query. (i.e., the output of the w column)

The canvas space will include an example output file. **Leave columns NULL/empty if the query does not give output for that column, or if the query is cyclic.**

Store this file together with your source code (in a folder `src/`) in a zip archive.

Name the file `<last-name-1>-<last-name-2>-impl.zip` or `<last-name>-impl.zip`.

IV. Evaluation

This assignment accounts for 40% of your final grade. There is no oral defense. Remember that you must hand in this assignment in order to pass for the course as a whole. **A late submission counts as 0/20 for the assignment.** The evaluation of your assignment takes into account the following aspects:

- **Requirements.** The correctness of your implementation, and the completeness of each requirement will be evaluated. Do not forget to include your output file, as part of the grading process is automated. We will, however, manually check your code as well.
- **Design.** The assignment allows for some creative liberty with respect to the specific implementation strategy of each requirement. Hence, make sure to put some thought into which data structures and algorithms you use, and make sure to motivate your decisions in comments, or in the report.
- **Code Quality.** As is typical in most programming assignments, the quality of your code is also an important factor. Pay specific attention to the clarity of your code, that is, use clear but concise variable and function names, and provide meaningful comments to motivate your design decisions.
Furthermore, make sure that your code is modularised and reusable. Use the abstraction facilities available in your language, and place care and attention into the prevention of needless code duplication.
- **Testing.** We expect that you provide some tests which showcase the functionality of your implementation in different scenarios.
- **Report.** The quality of your report will also be evaluated. Avoid paraphrasing or copy-pasting your code, as we have access to the source files; you can refer to your code instead. Make sure to focus on motivating your design.

Plagiarism

What is allowed? This assignment is to be made strictly individually and on your own. This means that you must create your assignment on an independent basis and you must be able to explain your work, reimplement it under supervision, and defend your solution. Copying work (code, text, etc.) from or sharing it with third parties (e.g., fellow students, websites, GitHub, etc.) is not allowed. Electronic tools will be used to compare all submissions with online resources and with each other, even across academic years.

Specifically, for this assignment:

- *It is allowed to use third-party code (e.g., logging libraries, testing libraries, etc.) with source attribution. If third-party code is used to fulfill any part of the requirements of the assignment (except the bonus feature on parsing CQs), **you will not receive any points for that part of the assignment.***
- *We do not consider your language's standard library to be third-party code. Hence, you may use that without attribution.*

What is plagiarism? Any action by a student that deviates from the instructions given and does not comply with the examination regulations is considered an irregularity. Plagiarism is also an irregularity. Plagiarism means the use of other people's work, adapted or otherwise, without careful acknowledgement of sources. (cf. OER, Article 118§2). Plagiarism may relate to various forms of works including text, code, images, etc.

What happens if plagiarism is suspected? Any suspicion of plagiarism will be reported to the Dean of Faculty without delay. Both user and provider of such code will be reported and will be dealt with according to the plagiarism rules of the examination regulations (cf. OER, Article 118). The dean may decide on (a combination of) the disciplinary sanctions, ranging from 0 out of 20 on the course to a prohibition from (re-)enrolling for one or multiple academic years (cf. OER, Article 118§5).

Contact us if you are in doubt as to whether or not something would be considered plagiarism.

Generative AI

Students are allowed to use generative AI tools to deepen their understanding of a topic, or to use it for parts of their assignment. However, **all content that was generated by such a tool should be included in an appendix to your report**, and it should be **cited appropriately** throughout your assignment. Our primary objective is to determine whether *you, the student*, understands the material; we do not wish to test whether the AI tools' training datasets are up-to-date.

Concretely, per piece of content that was generated and copied into your code or report, you have to include the following information:

- Which **tool** was used to generate this content? (e.g., ChatGPT, DALL-E, Copilot, ...)
- What was the **prompt** given to this generative AI tool?
- What was the **output** given by this generative AI tool?
- A way of uniquely identifying the prompt in citations (e.g., prompt#1, or IEEE citation numbers such as [4])

You also have to include a general declaration that you made use of generative AI in your assignment, and you have to explain for which purpose; you can find an example below. Failure to provide this information, and failure to cite it, will be treated like an irregularity (and thus, plagiarism) under the VUB's teaching and exam regulations.

Example Appendix Entry

For this assignment I made use of “ChatGPT-3.5” (<https://openai.com/chatgpt>), for [purpose of the use]. In my text, I indicated where I used the generated content using the following citation style: ([name-of-ai-tool], prompt#1). Which refers to the prompts given below.

Tool: ChatGPT-3.5

Prompt Number: 1

Prompt: *“Can you define what a semi-join is?”*

Output:

“Certainly! A semi-join is a type of database join operation that filters the rows from the first table based on the existence of matching rows in the second table, but it does not include the actual data from the second table in the result set.”

Tool: ChatGPT-3.5

Prompt Number: 2

Prompt: *“Can you define what an outer join is?”*

Output:

“An outer join is a type of database join operation that retrieves records from one or more tables, even if there is no match between the columns being joined. Unlike inner joins, which only return rows with matching values in both tables, outer joins include unmatched rows from one or both of the tables in the result set.”

Example Citation

APA-style: (ChatGPT, prompt#1)

IEEE-style: [1] where the prompt entry in the appendix is preceded by [1].