

csc343 winter 2020

assignment #1: relational algebra

sample solutions

goals

This assignment aims to help you learn to:

- read a relational scheme and analyze instances of the schema
- read and apply integrity constraints
- express queries and integrity constraints of your own
- think about the limits of what can be expressed in relational algebra

Your assignment must be typed to produce a PDF document **a1.pdf** (hand-written submissions are not acceptable). You may work on the assignment in groups of 1 or 2, and submit a single assignment for the entire group on [MarkUs](#). You must establish your group well before the due date by submitting an incomplete, or even empty, submission.

background

You will be working on a schema and queries for a database used by a zoological institute to track an archive of their artifacts.

During a field trip collectors gather a variety of artifacts of the animals they study, resulting in tissue samples, images, physical models (such as casts of paw prints), or live colonies.

After arriving at the institute, artifacts must be safely stored and maintained by technicians. Some artifacts are cited in one or more publications. In all cases the official species name must be recorded, and must appear in the [Catalogue of Life database](#). If correct taxonomic practices are followed, each species belongs to exactly one genus, and each genus to exactly one family. Tables COL, Genus, and Species are derived from Catalogue of Life database.

relations

- Collection(CID, date, SID)
Tuples here represent entire collections from a field trip, where *CID* is the collection ID, *date* is the starting date of the field trip, and *SID* is the staff ID of the collector.

- **Collected**(CID, AN)
A tuple here represents the fact that collection *CID* includes artifact number *AN*. A single collection usually contains multiple artifacts, and a single artifact may be aggregated from more than one collection.
- **Artifact**(AN, species, type, location, SID)
Tuples here represent single artifact collected in the field. *AN* is the artifact number, *species* is the scientific species name, *type* is one of tissue, image, model, or live, *location* is where it was collected, and *SID* is the staff number of the technician who maintains this artifact.
- **Published**(AN, journal, date)
A tuple here represents the fact that artifact *AN* was mentioned in scholarly publication *journal* with publication date *date*.
- **Staff**(SID, name, email, rank, date)
These tuples represent a member of the institute's scientific staff. *SID* is the staff ID, *name* is their full name, *email* is their professional email, *rank* is one of: technician, student, pre-tenure, or tenured, and *date* is the date when they attained that rank.
- **COL**(family)
A singleton tuple here means that *family* is a scientific zoological family name that appears in the Catalogue of Life.
- **Genus**(genus, family)
A tuple here means that *genus* is in family *family*.
- **Species**(species, genus)
A tuple here means that *species* is in genus *genus*.

our constraints

For each of the following constraints give a one sentence explanation of what the constraint implies, and why it is required.

- $\Pi_{species}(Artifact) - \Pi_{species}(Species) = \emptyset$.
sample solution: Every species maintained in Artifact is in the Species table, to ensure that standard taxonomic names are used.
- $\Pi_{rank}(Staff) \subseteq \{'technician', 'student', 'pre-tenure', 'tenure'\}$.
sample solution: All staff in Staff have one of these four titles, to ensure that properly qualified personnel work on collections.
- $\Pi_{family}(Genus) - \Pi_{family}(COL) = \emptyset$.
sample solution: Every family in table Genus is in the Catalogue of Life, to ensure standard taxonomic names are used.
- $\Pi_{genus}(Species) \subseteq \Pi_{genus}(Genus)$.
sample solution: Every genus in Species is in table Genus, to ensure standard taxonomic names are used.

- $\Pi_{CID}(Collected) = \Pi_{CID}(Collection)$.

sample solution: Every collection/artifact pair refers to a collection, to ensure that irrelevant pairs are not recorded.

- $\Pi_{AN}(Artifact) = \Pi_{AN}(Collected)$.

sample solution: Every artifact is from at least one collection, to ensure that no under-documented artifacts are present.

- $\Pi_{SID}(Collection) \subseteq \Pi_{SID}(Staff)$.

sample solution: All collectors are on staff, to ensure that no informal, free-lance collectors are at large.

- $\Pi_{SID}(Artifact) \subseteq \Pi_{SID}(Staff)$.

sample solution: All maintainers are on staff, to ensure that no free-lance maintainers are working.

- $\Pi_{type}(Artifact) \subseteq \{'tissue', 'image', 'model', 'live'\}$

sample solution: Each artifact is of one of these four types, to limit the complexity of types of data recorded.

- $\Pi_{AN}(Published) \subseteq \Pi_{AN}(Artifact)$

sample solution: All published artifacts are recorded in Artifact, to ensure that anything published has been properly documented.

queries

Write relational algebra expressions for each of the queries below. You must use notations from this course and operators:

$$\Pi, \sigma, \rho, \bowtie, \bowtie_{condition}, \times, \cap, \cup, -, =$$

You may also use constants:

$$\text{today (for current date)} \quad \emptyset \text{ (for the empty set)}$$

In your queries pay attention to the following:

- All relations are sets, and you may only use relational algebra operators covered in Chapter 2 of the course text.
- Do not make assumptions that are not enforced by our constraints above, so your queries should work correctly for any database that obeys our schema and constraints.
- Other than constants such as 23 or "lupus", a select operation only examines values contained in a tuple, not aggregated over an entire column.
- Your selection conditions can use arithmetic operators, such as $+$, \leq , \neq , \geq , $>$, $<$ and friends. You can use logical operators such as \vee , \wedge , and \neg , and treat dates and numeric attributes as numbers that you can perform arithmetic on.
- Use good variable names and provide lots of comments to explain your intentions.

- Return multiple tuples if that is appropriate for your query.

There may be a query or queries that cannot be expressed in the relational algebra you have been taught so far, in which case just write “cannot be expressed.” The queries below are not in any particular order.

1. Rationale: Performance reviews include seeing how current the work is of staff who have held their current rank for a long time.

Query: Find the most recent collection date of any artifact collected by a staff member who has held their current rank the longest. Keep ties.

sample solution:

```
// how long rank has been held:
rank-duration(SID, duration) :=  $\Pi_{SID, today-date}(\text{Staff})$ 

// not longest-serving staff
not-longest(SID) :=  $\Pi_{SID} [\sigma_{r1.duration < r2.duration} (\rho_{r1} \text{rank-duration} \times \rho_{r2} \text{rank-duration})]$ 

// longest-serving staff:
longest(SID) :=  $\Pi_{SID} (\text{rank-duration} - \text{not-longest})$ 

// collection dates for longest-serving staff:
collection-dates(SID, date) :=  $\Pi_{SID, date} (\text{longest} \bowtie \text{Collection})$ 

// not-newest collection:
not-newest(SID, date)
:=  $\Pi_{s1.SID, s1.date} [\sigma_{s1.SID = s2.SID \wedge s1.date < s2.date} (\rho_{s1} \text{collection-dates} \times \rho_{s2} \text{collection-dates})]$ 

// newest collection:
answer(SID, date) := collection-dates - not-newest
answer
```

2. Rationale: Staff who maintain every artifact in some collection should be considered favourably in performance reviews.

Query: Find all staff who maintain all artifacts in at least one collection.

sample solution:

```
// maintainers of collections
MaintainCollect :=  $\Pi_{CID, SID} (\text{Collected} \bowtie \text{Artifact})$ 

// multiple maintainer collections
MultipleMaintainCollect
:=  $(\rho_{mc1} \text{MaintainCollect}) \bowtie_{mc1.SID \neq mc2.SID \wedge mc1.CID = mc2.CID} (\rho_{mc2} \text{MaintainCollect})$ 

// single maintainer collections
SingleMaintainCollect :=  $(\Pi_{CID} \text{MaintainCollect}) - (\Pi_{CID} \text{MultipleMaintainCollect})$ 

// staff who maintain entire collections
 $\Pi_{SID} (\text{SingleMaintainCollect} \bowtie \text{Collected} \bowtie \text{Artifact})$ 
```

3. Rationale: An artifact collected and maintained by the same staff may have some special requirements that should be investigated.

Query: Find all artifacts that were collected by the same staff who maintains them.

sample solution:

// artifacts where staff in Collection and Artifact are the same:

$$\Pi_{AN} (\text{Collection} \bowtie \text{Collected} \bowtie \text{Artifact})$$

4. Rationale: Identify multi-talented field workers.

Query: Find all staff who have collected at least 3 artifacts from every species in some family.

sample solution:

// artifact info:

artifact-species(AN, species, SID)

$$:= \Pi_{AN, species, Collection.SID} (\text{Artifact} \bowtie \text{Collected} \bowtie_{\text{Collected.CID}=\text{Collection.CID}} \text{Collection})$$

// at least two of a species:

at-least-two(AN, species, SID)

$$:= \Pi_{a1.AN, a1.species, a1.SID} [\sigma_{a1.AN < a2.AN \wedge a1.species = a2.species \wedge a1.SID = a2.SID} (\rho_{a1} \text{artifact-species} \times \rho_{a2} \text{artifact-species})]$$

// at least three of a species:

at-least-three(SID, species)

$$:= \Pi_{a1.SID, a1.species} [\sigma_{a1.AN < a2.AN \wedge a1.species = a2.species \wedge a1.SID = a2.SID} (\rho_{a2} \text{at-least-two} \times \rho_{a1} \text{artifact-species})]$$

// family for at-least-three:

$$\text{real}(\text{SID}, \text{species}, \text{family}) := (\text{Species} \bowtie \text{Genus} \bowtie \text{at-least-three})$$

// ideal case:

$$\text{ideal}(\text{SID}, \text{species}, \text{family}) := \Pi_{SID, species, family} [(\Pi_{SID, family} \text{real}) \bowtie \text{Genus} \bowtie \text{Species}]$$

// families missing at least one species:

$$\text{missing}(\text{SID}, \text{family}) := \Pi_{SID, family} [\text{ideal} - \text{real}]$$

// staff with complete family:

$$\Pi_{SID} [\Pi_{SID, family} (\text{real}) - \text{missing}]$$

5. Rationale: Which publications might have some specialized niche focus?

Query: Find all publications that have used exactly 2 of our artifacts.

sample solution:

```
// publications that used at least 3 artifacts

TripleProduct :=  $\rho_{p1}(\text{Published}) \times \rho_{p2}(\text{Published}) \times \rho_{p3}(\text{Published})$ 
AtLeast3 :=  $\sigma_{p1.AN < p2.AN < p3.AN \wedge p1.journal = p2.journal = p3.journal}$  TripleProduct

// publications that used at least 2 artifacts

AtLeast2 :=  $\sigma_{p1.AN < p2.AN \wedge p1.journal = p2.journal}$  TripleProduct

// publications that used exactly 2 artifacts

 $\Pi_{p1.journal}(\text{AtLeast2}) - \Pi_{p1.journal}(\text{AtLeast3})$ 
```

6. Rationale: Identify motherlode locations.

Query: Find all locations where at least one artifact from every family has been collected.

sample solution:

```
// all actual location/family pairs

actual-location-family :=  $\Pi_{location, family}(\text{Artifact} \bowtie \text{Species} \bowtie \text{Genus})$ 

// all possible location/family pairs

possible-location-family :=  $\Pi_{location, family}(\text{Artifact} \times \text{COL})$ 

// missed location/family pairs

missed-location-family := possible-location-family - actual-location-family

// locations with every family

 $(\Pi_{location} \text{actual-location-family}) - (\Pi_{location} \text{missed-location-family})$ 
```

7. Rationale: Exclusively tissue sample collectors may need extra support for special reagents and shipping costs.

Query: Find all staff who have collected only tissue samples.

sample solution:

```
// staff who have collected non-tissue

NonTissueCollectorSID :=  $\Pi_{Collection.SID}((\sigma_{type \neq tissue} \text{Artifact}) \bowtie \text{Collected} \bowtie \text{Collection})$ 

// all collectors

CollectorSID :=  $\Pi_{SID} \text{Collection}$ 

// collectors who never collected non-tissue

CollectorSID - NonTissueCollectorSID
```

8. Rationale: Collection staff who should be encouraged to diversify their network.

Query: Find all staff pairs who have worked only with each other on collections.

sample solution: We are going to follow the assumption that pairs of collector-maintainer or maintainer-collector who only work with each other, and neither works with themselves. We will also accept solutions that explicitly state that they are considering maintainer/maintainer collaborations

```
// collector-maintainer pairs

CollectMaintain( $S1, S2$ ) :=  $\Pi_{C.SID, A.SID}(\sigma_{C.SID < A.SID}(\rho_C \text{Collection} \bowtie \text{Collected}) \times \rho_A(\text{Artifact}))$ 

// maintainer-collector pairs

MaintainCollect( $S1, S2$ ) :=  $\Pi_{A.SID, C.SID}(\sigma_{C.SID > A.SID}(\rho_C \text{Collection} \bowtie \text{Collected}) \times \rho_A(\text{Artifact}))$ 

// all pairs

AllPairs := CollectMaintain  $\cup$  MaintainCollect

// product of all pairs

PairProduct :=  $(\rho_A \text{AllPairs}) \times (\rho_B \text{AllPairs})$ 

// non-exclusive pairs

NonExclusive1 :=  $\sigma_{A.S1=B.S2 \vee A.S2=B.S1} \text{PairProduct}$ 
NonExclusive2 :=  $\sigma_{A.S1=B.S1 \wedge A.S2 \neq B.S2} \text{PairProduct}$ 
NonExclusive3 :=  $\sigma_{A.S2=B.S2 \wedge A.S1 \neq B.S1} \text{PairProduct}$ 
NonExclusive := NonExclusive1  $\cup$  NonExclusive2  $\cup$  NonExclusive3

// exclusive pairs

AllPairs -  $\Pi_{S1, S2}(\rho_{S1 \leftarrow A.S1, S2 \leftarrow A.S2} \text{NonExclusive})$ 
```

9. Rationale: Track the influence of a given staff member.

Query: Staff member SID_1 is influenced by staff member SID_2 if (a) they have ever worked together on a collection or (b) if SID_1 has ever worked with a staff member who is influenced by SID_2 . Find $SIDs$ of staff members influenced by SID_2 .

sample solution: This query cannot be expressed in our RA. This would require an arbitrary number of joins or cartesian products to find the transitive closure of the “influenced” relation.

your constraints

For each of these constraints you should derive a relational algebra expression of the form $R = \emptyset$, where R may be derived in several steps, by assigning intermediate results to a variable. If the constraint cannot be expressed in the relational algebra you have been taught, write “cannot be expressed.”

1. No species is also a genus.

sample solution:

$$\Pi_{species}(\text{Species}) \cap \Pi_{genus}(\text{Genus}) = \emptyset$$

2. No genus belongs to more than one family.

sample solution:

$$(\rho_{g1} \text{Genus}) \bowtie_{g1.family \neq g2.family \wedge g1.genus = g2.genus} (\rho_{g2} \text{Genus}) = \emptyset$$

3. All publications must be published after all artifacts they use have been collected.

sample solution:

// collection dates for every artifact in publication:

```

dates(collection-date, publication-date)
:=  $\Pi_{c1.date, p.date} (\rho_{c1} \text{Collection} \bowtie_{c1.CID=c2.CID} \rho_{c2} \text{Collected} \bowtie_{c2.AN=p.AN} \rho_p \text{Publication})$ 
 $\sigma_{\text{collection-date} > \text{publication-date}} \text{dates} = \emptyset$ 

```

4. Students may not maintain live artifacts.

sample solution:

$$\sigma_{\text{type}='live' \wedge a.SID=s.SID \wedge s.rank='student'} [\rho_a(\text{Artifact}) \times \rho_s(\text{Staff})] = \emptyset$$

submissions

Submit **a1.pdf** on **MarkUs**. One submission per group, whether a group is one or two people. You declare a group by submitting an empty, or partial, file, and this should be done well before the due date. You may always replace such a file with a better version, until the due date.

Double check that you have submitted the correct version of your file by downloading it from MarkUs.

marking

We mark your submission for correctness, but also for good form:

- For full marks you should add comments to describe the *data*, rather than *technique*, of your queries. These may help you get part marks if there is a flaw in your query.
- Please use the assignment operator, “:=” for intermediate results.
- Name relations and attributes in a manner that helps the reader remember their intended meaning.
- Format the algebraic expressions with line breaks and formatting that help make the meaning clear.