Relational Calculus & Relational Algebra

- 1. Give one reason why MRN is a poor choice for the primary & foreign keys? One single patient may come to see clinicians for more than one time, so the same MRN could appear on multiple tuples in the relation VISIT.
- 2. Write Relational Calculus expressions for the following:
 - (a) Who is 25 years old?

```
{pa. FIRSTNAME, pa. LASTNAME | PATIENT(pa) \land pa. AGE = '25'}
```

(b) Who had a medical visit in December, 2017?

```
{ pa. FIRSTNAME, pa. LASTNAME | PATIENT(pa) \land \exists (v)(VISIT(v) \land v. DATETIME = 'December, 2017' \land pa. MRN = v. MRN)}
```

(c) Who got a 'flu shot'?

```
{pa. FIRSTNAME, pa. LASTNAME | PATIENT(pa) \land \exists (v, pr)(VISIT(v) \land PROCEDURE(pr) \land pr. NAME = 'flu shot' <math>\land v. VISIT\_ID = pr. VISIT\_ID \land pa. MRN = v. MRN)}
```

(d) Who did NOT get a 'flu shot'?

```
{pa. FIRSTNAME, pa. LASTNAME | PATIENT(pa) \land \neg \exists (v, pr)(VISIT(v) \land PROCEDURE(pr) \land pr. NAME = 'flu shot' <math>\land v. VISIT\_ID = pr. VISIT\_ID \land pa. MRN = v. MRN)}
```

(e) What is the first and last name of all patients who have seen MD Paula Jones?

```
{pa. FIRSTNAME, pa. LASTNAME | PATIENT(pa) \land \exists (v, pr, c)(VISIT(v) \land PROCEDURE(pr) \land CLINICIAN(c) \land c. CERT = 'MD' \land c. FIRSTNAME = 'Paula' <math>\land c. LASTNAME = 'Jones' \land pr. CLIN_ID = c. CLIN_ID \land v. VISIT_ID = pr. VISIT_ID \land pa. MRN = v. MRN)}
```

- 3. Write Relational Algebra expressions for the following questions. Return the relation primary key to identify the tuples, unless otherwise specified.
 - (a) Who is 25 years old?

```
\pi_{MRN}(\sigma_{AGE} = `25", (PATIENT))
```

(b) What is the first and last name of all patients who have seen a Physician's Assistant (cert is 'PA')

```
C \leftarrow \rho firstname c/firstname, lastname c/lastname(CLINICIAN)
```

```
TIFIRSTNAME, LASTNAME (OCERT = 'PA' (PATIENT * VISIT * PROCEDURE * C))
```

(c) Which patients have the same names and age?

$$\Pi$$
 P.MRN, P1.MRN(σ P.MRN != P1.MRN (ρ P(...)(Patient) \bowtie P.Firstname=P1.Firstname \land P.Lastname=P1.Lastname \land P.Age=P1.Age ρ P1(....)(Patient))

(d) Which patients who got a flu shot also got a measles immunization during the same visit?

$$R \leftarrow \pi_{\text{VISIT_ID, NAME1}}(\rho_{\text{NAME1/NAME}}(\sigma_{\text{NAME}} = \text{`flu shot'}(PROCEDURE)))$$

$$S \leftarrow \pi_{VISIT_ID, NAME2}(\rho_{NAME2/NAME}(\sigma_{NAME} = 'measles immunization', (PROCEDURE)))$$

 $\pi_{MRN}(R*S*VISIT*PATIENT)$

(e) Which patients who have seen an MD have not seen a PA?

$$C \leftarrow \rho_{FIRSTNAME_C/FIRSTNAME, LASTNAME_C/LASTNAME}(CLINICIAN)$$

$$R \leftarrow \pi_{FIRSTNAME, LASTNAME} (\sigma_{CERT = 'MD'}(PATIENT * VISIT * PROCEDURE * C))$$

$$\pi_{MRN}$$
 (($\pi_{FIRSTNAME}$, Lastname(PATIENT) — s) $rac{R}$

Queries

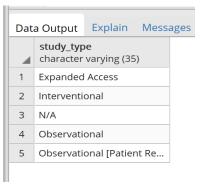
1. List the nct id and study type from the study whose brief title is "Autologous Cell Therapy After Stroke".

SELECT nct_id, study_type
FROM studies
WHERE brief title = 'Autologous Cell Therapy After Stroke';

Data Output		Explain	Messages	Notifications
4	nct_id character (11)		study_type character varying (35)	
1	NCT00908856		Interventional	
		,		

2. List the different values for study type, in alphabetical order.

SELECT DISTINCT study_type FROM studies ORDER BY study_type ASC;



3. How many terminated studies that started and completed in 2016 have reported events?

SELECT COUNT (DISTINCT studies.nct_id)

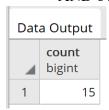
FROM studies, reported events

WHERE studies.nct_id = reported_events.nct_id

AND start date >= '2016-01-01'

AND completion date < '2017-01-01'

AND overall status = 'Terminated';



4. How many of the studies that started in February 2016, but on or after the 15th, are expected to complete (or have completed) within 6 months of their start date?

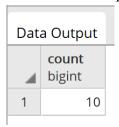
SELECT COUNT(DISTINCT nct_id)

FROM studies

WHERE start date >= '2016-02-15'

AND start date < '2016-03-01'

AND completion date <= start date + INTERVAL '6 MONTH';



Reading/Short answer

- 1. List three reasons the database was "normalized."
 - (1) Normalization makes the organization of data more efficient.
 - (2) Reduce redundancy.
 - (3) Ensure logical data dependencies by storing only related data within a given table.
- 2. On page 4, the authors list a number of data elements that are contained in the XML download files from ClinicalTrials.gov. How are these different elements implemented in the design table in the AACT database?

The values associated with each data element are tested for correctness and completeness by comparing them with the original source data from downloaded XML files.

3. Which non-description field in the design table is least populated? That is, which field is most often left blank?

observational model

- 4. In the conditions table there is an attribute called "name" and an attribute called "downcase name". What is the difference? Why might a database provide both of these fields? What trade-offs are involved?
- (1) For the attribute "name", the first letter of the contents is uppercase or the contents are uppercase acronyms. While in "downcase name", all letters are lowercase.
- (2) It's more convenient for a user to query on database.
- (3) It requires more space to store the information.
- 5. Look at some of the name, downcase name pairings. Do you see any anomalies? Give 2 examples of what challenges might these anomalies pose to a user of the database.
- (1) Some values of the "name" do not have pairing in "downcase name". If users want to run some queries using "downcase name", they might miss some of the data they want. For example, from row 29 to 36, the values of their "downcase name" is null. But all of these rows have their corresponding value in attributes "name".
- (2) The ambiguity of name between uppercase and lowercase. For example, the attribute "name" of row 2 is "Hiv". However, the name of row 19 is "HIV". Actually both of them refer to the same disease. So a user looking for the study of this disease might miss one of the two HIV.