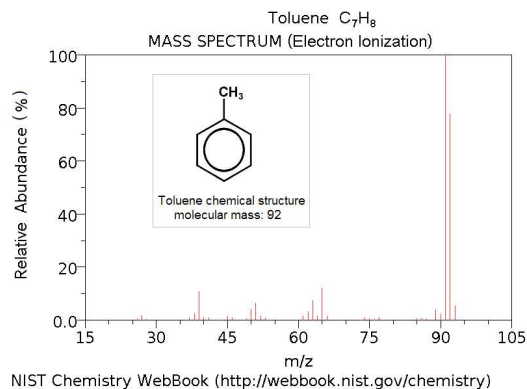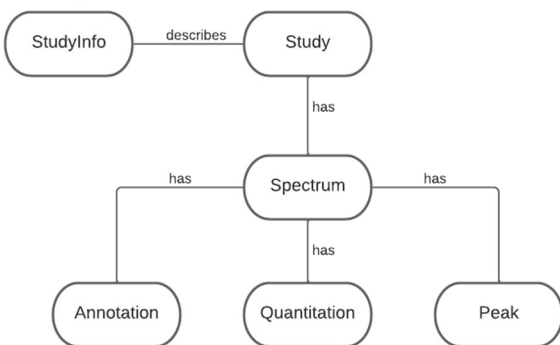# Lab 9

During this and the next few labs, we'll create a database to store metabolomics mass-spectrometry data.

## From Wikipedia:

- **Metabolomics** is the scientific study of chemical processes involving metabolites, the small molecule substrates, intermediates and products of cell metabolism.
- **Mass spectrometry** (MS) is an analytical technique that is used to measure the mass-to-charge ratio of ions. In a typical MS procedure, a sample is ionized by bombarding it with a beam of electrons. This causes some of the sample's molecules to break up into positively or negatively charged fragments. These ions (fragments) are then separated according to their mass-to-charge ratio (m/z). The results are presented as a mass spectrum, a plot of intensity as a function of the mass-to-charge ratio (m/z).



Resulting mass spectra are (mostly) unique for each molecule, and can be used for identifying/annotating unknown molecules.

## Conceptual schema of the database:



## Lab 9 task:

In this lab, we'll focus on three entities: **Study**, **Spectrum**, and **Peak**. They are related as follows:

- Each study has multiple spectra. Each spectrum belongs to a single study.
- Each spectrum has multiple peaks. Each peak belongs to a single spectrum.

# Task 1.

Download and unzip the data archive from Canvas. It contains 4 folders, each folder contains data from a single study. In each folder, you'll see the following files:

- **meta.json** contains information (title, summary, url) about each study,

- **spectra.msp** contains some information about each spectrum and its peaks (m/z-intensity pairs)

You can open these files in any text editor, but I recommend Visual Studio Code (https://code.visualstudio.com/). Take a look at these files and make sure you understand their content.

## Task 2 (15pts).

Design tables for the entities **Study**, **Spectrum**, and **Peak**. Draw an Entity-Relationship diagram (ERD). In ERD, specify name of each table, list of attributes, primary and foreign keys, and relationships between the tables.
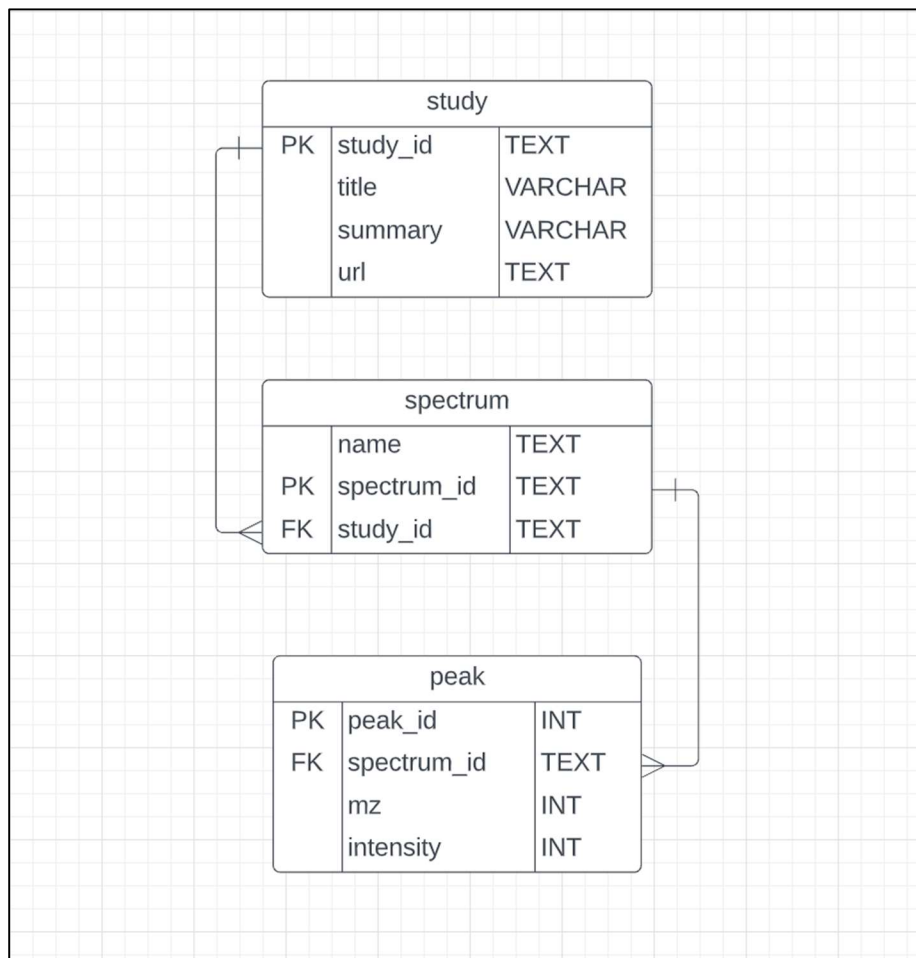


Fig: ERD for study, spectrum and peak

It is up to you to decide what attributes you want to have in each table.

- At the minimum, you need the attributes to be able to execute the SQL queries from Task 4,
- I'd also suggest to add attributes for study title, study summary, study ID, and spectrum name,
- The **Peak** entity should have attributes **mz** and **intensity**.

Course: BINF-6211 – Lab 9
Name: Drashti Mehta


Report the Entity-Relationship diagram.

Here, I'm creating a temporary table for importing unique data in the Spectrum and Peak tables. This "temp_table" has information from all the spectra files (created as shown in part 3) and later data is imported in the tables. The way it works is that the data from every file is imported into this temp_table one by one and it is taken into Spectrum and Peak table. Study_id can be added manually as default as and when each file is imported.

```
CREATE TABLE Study (
study_id CHAR PRIMARY KEY,
title VARCHAR,
summary VARCHAR,
url TEXT
);

CREATE TABLE Spectrum (
name TEXT,
spectrum_id TEXT PRIMARY KEY,
study_id TEXT REFERENCES Study (study_id)
);

CREATE TABLE temp_table (
name TEXT,
id TEXT PRIMARY KEY,
num_peaks INTEGER,
mz INTEGER,
intensity INTEGER
);

CREATE TABLE Peak (
peak_id INTEGER PRIMARY KEY AUTOINCREMENT,
spectrum_id TEXT REFERENCES Spectrum (spectrum_id),
mz INTEGER,
intensity INTEGER
);
```


# Task 3 (15pts).

Fill out the tables that you created with data from the archive. In order to read the .msp files, I've attached scripts both in Python and R to convert those files into a csv.

You can run the Python script in the terminal (you may need to install Pandas):
> python3 read_msp.py PATH_TO_MSP

Course: BINF-6211 – Lab 9
Name: Drashti Mehta

You can run the R script in RStudio or in the terminal (you need to edit paths to the files inside the R script):
> Rscript read_msp.py

Upload the database to Canvas.

ANSWER:
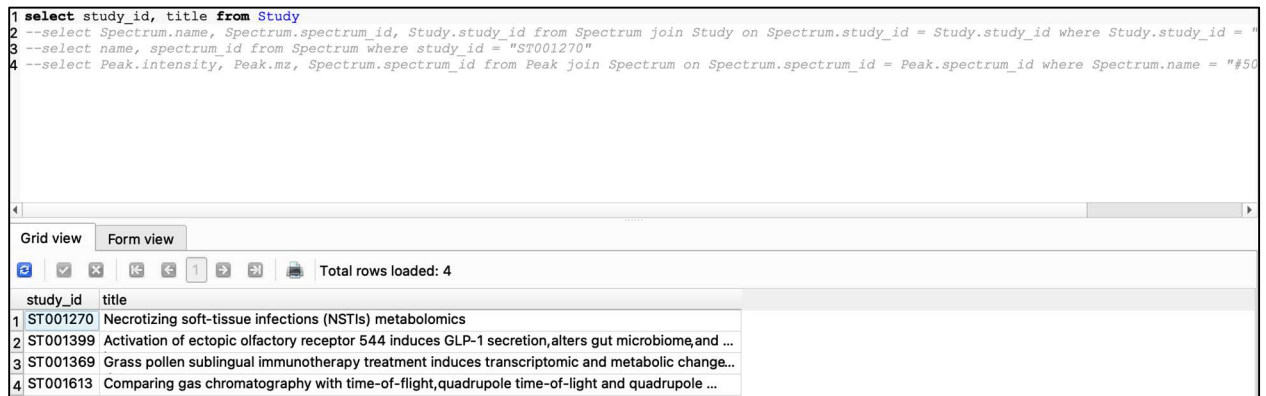The files are attached to the submission section with the database.

# Task 4 (15pts).
Try to execute the following SQL queries:
- List all studies in the database
- List all spectra for the study "ST001270"
- List all peaks for the spectrum "#50 75.0000@6.0753 MS1+"

Report the output of these queries. You can take a screenshot of the first few rows of the output.

ANSWER:

1. select study_id, title from Study

```
1 select study_id, title from Study
2 --select Spectrum.name, Spectrum.spectrum_id, Study.study_id from Spectrum join Study on Spectrum.study_id = Study.study_id where Study.study_id = "
3 --select name, spectrum_id from Spectrum where study_id = "ST001270"
4 --select Peak.intensity, Peak.mz, Spectrum.spectrum_id from Peak join Spectrum on Spectrum.spectrum_id = Peak.spectrum_id where Spectrum.name = "#50
```

Grid view | Form view

Total rows loaded: 4

| | study_id | title |
|---|---|---|
| 1 | ST001270 | Necrotizing soft-tissue infections (NSTIs) metabolomics |
| 2 | ST001399 | Activation of ectopic olfactory receptor 544 induces GLP-1 secretion,alters gut microbiome,and ... |
| 3 | ST001369 | Grass pollen sublingual immunotherapy treatment induces transcriptomic and metabolic change... |
| 4 | ST001613 | Comparing gas chromatography with time-of-flight,quadrupole time-of-light and quadrupole ... |

2. This can be done in two ways but way 1 is preferred.

1) select Spectrum.name, Spectrum.spectrum_id, Study.study_id from Spectrum join Study on Spectrum.study_id = Study.study_id where Study.study_id = "ST001270"

OR

2) select name, spectrum_id from Spectrum where study_id = "ST001270"

3. select Peak.intensity, Peak.mz, Spectrum.spectrum_id from Peak join Spectrum on Spectrum.spectrum_id = Peak.spectrum_id where Spectrum.name = "#50 75.0000@6.0753 MS1+"