

# Foundations of Data Science for Biologists

## More SQL

BIO 724D

2024-NOV-19

Instructors: Greg Wray and Paul Magwene

# Creating and maintaining an SQL database

# SQL is limited to a small set of operations

Only four kinds of operations are permitted with relational databases

Create, Read, Update, Delete (CRUD)

**CREATE** to create a new database, table, or relation; **INSERT** to add rows to a table

**SELECT** to query (retrieve information from) a database ✓

**UPDATE** to change information in a table

**DROP** to remove a database, table, or relation; **DELETE** to remove rows from a table

SQL has some additional keywords, but most work with the above set

Clauses within statements: **LIMIT**, **WHERE**, **JOIN**, **GROUP BY**, etc.

Operators within statements: **=**, **>**, **AND**, **NOT**, **LIKE**, **BETWEEN**, etc.

Functions within statements: **MIN**, **MEAN**, **COUNT**, **DISTINCT**, etc.

**dplyr** and **Pandas** implement some SQL-like functions (and were inspired by it)

# Constraints

# Constraints are designed to ensure data integrity

**Constraints** are rules that govern what data can be present and how data interact

SQL provides several mechanisms that are easy to implement but powerful

Constraints should be defined when you create a table

Some SQL implementations allow constraints to be added later

But SQL will complain if the existing data don't conform to the new rule!

# Column (single-table) constraints

**UNIQUE**: ensures that every value in a column is distinct from all others

**NOT NULL**: ensures that a column contains no NULL values

**PRIMARY KEY**: both **UNIQUE** and **NOT NULL**; also important for **normalization** (later)

**DEFAULT**: inserts a default value if no value is specified during **INSERT** or **UPDATE** operations

**CHECK**: ensures that the value in a column meets a specified condition(s)

# Two-table constraints

**FOREIGN KEY:** prevents actions that would destroy a **relation** between two tables

a foreign key depends on the values  
in a column in another table

the referenced column in the other  
table must contain unique values

observations

seq	genus	species	date	location
1023	Pycnonotus	tricolor	2007-06-14	Fairview Hotel
1024	Milvus	migrans	2007-06-14	Fairview Hotel
1032	Chalcomitra	amethystina	2007-06-14	Sheldrick Centre
1033	Pycnonotus	tricolor	2007-06-14	Sheldrick Centre

locations

location	prov	country	clim	elev	geolocation
Everard Reserve	Vic	Australia	Csb	90	-37.68,145.49
Everglades NP	FL	USA	Aw	2	25.39,-80.63
Fairview Hotel	Nb	Kenya	Cfb	1715	-1.29,36.80
Faskrudsfjordur	Au	Iceland	ET	20	64.93,-14.01

A foreign key prevents:

Adding an observation with a location that does not exist in the locations table

Removing a location from the locations table that is referenced by an observation

# NULL is a special value in SQL

NULL means the value does not exist in the database

Think of it as meaning “unknown” or “missing information”

It does *not* mean 0 or FALSE

In most cases, operations with NULL values “propagate”

1 + NULL returns NULL

NULL / 0 most SQL implementations return NULL

'a' || NULL || 'b' returns NULL

Technically, NULL is a third logic value in SQL

This is related to the algebra of joins and is beyond the scope of our class



# Database design and normalization

# Database design

A database should be organized to:

- Minimize mistakes when entering and updating data

- Minimize redundant information

- Simplify the query process

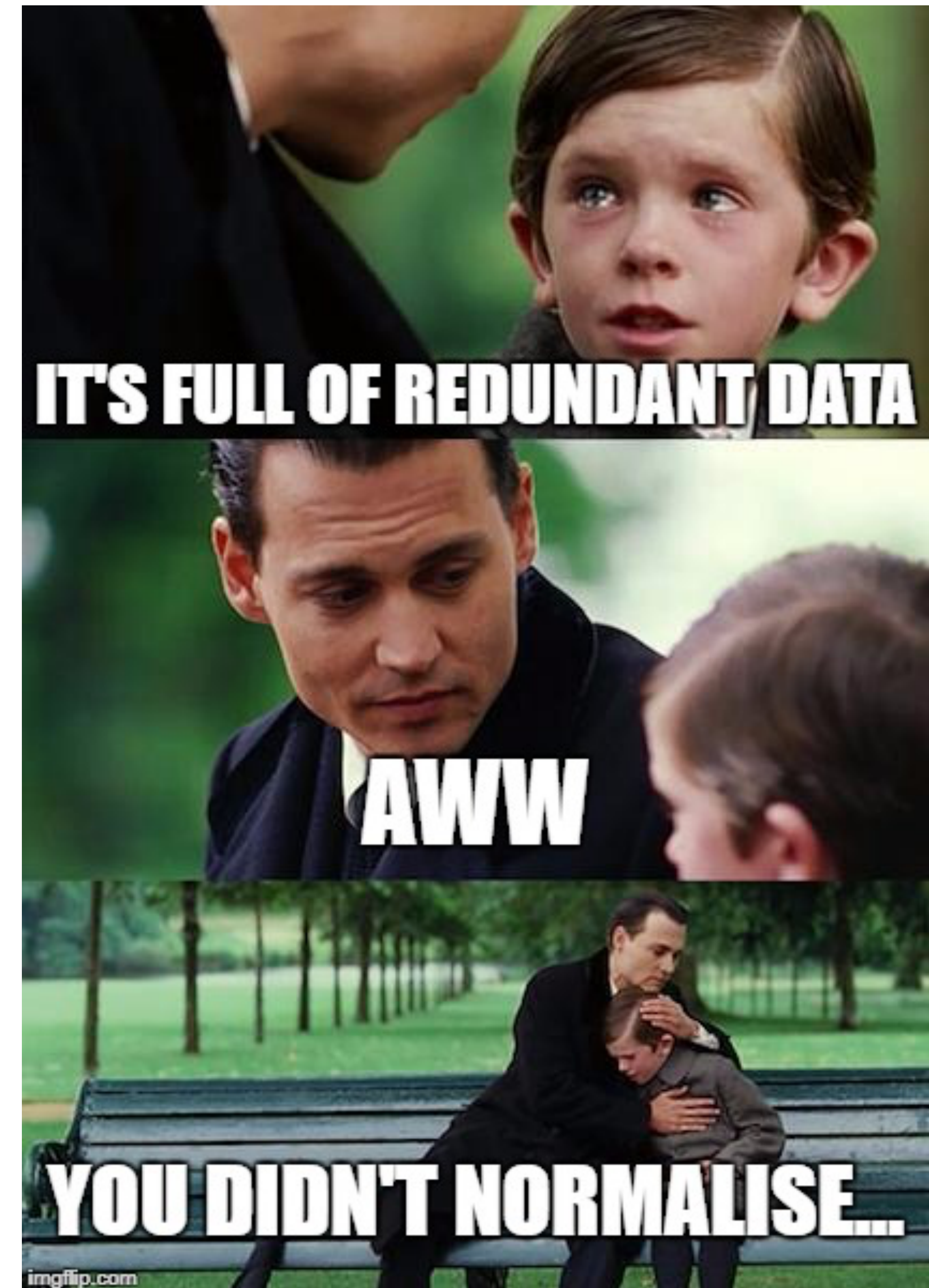
The solution to all these goals: database **normalization**

Normalization helps you decide:

- How many tables?

- What should be their primary keys?

- What relations should link the tables?



# First normal form

## Do not use row order to convey information

Row order is neither stable nor predictable

If consistent row order is needed for output, include an index column

## Do not mix data types within a column

Some spreadsheet programs allow this (looking at you, Excel)

As a result, you may encounter errors when importing data

## Do not use repeating groups

Examples of repeating groups: vectors, lists, tuples, and dictionaries

These values need to go into another table

## A primary key is required

Values must be unique within a single column or together across multiple columns

The most natural way to think about a PK is as a row identifier



# Second and third normal forms

## Second normal form

Concerns how non-PK columns relate to the PK

Each non-PK must depend on the PK (i.e., its value is not a function of something else)

Each non-PK must depend on the entire PK (i.e., *both* columns if a multi-column PK)

If this is not the case, the column in question should be in another table

## Third normal form

Concerns transitive dependencies: non-PK column depends on another non-PK column

Each non-PK must depend only on the PK

If this is not the case, the column in question should be in another table

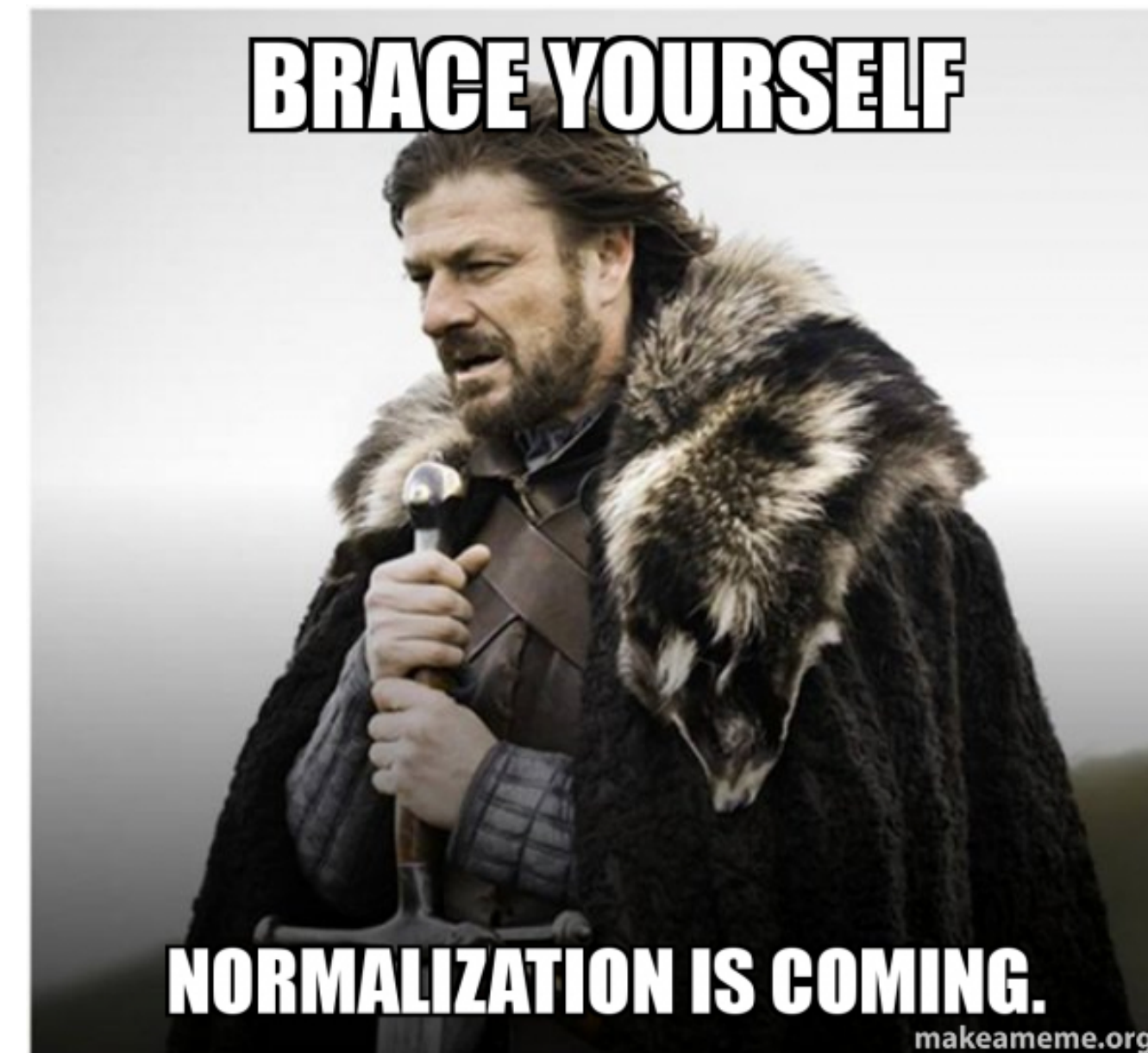
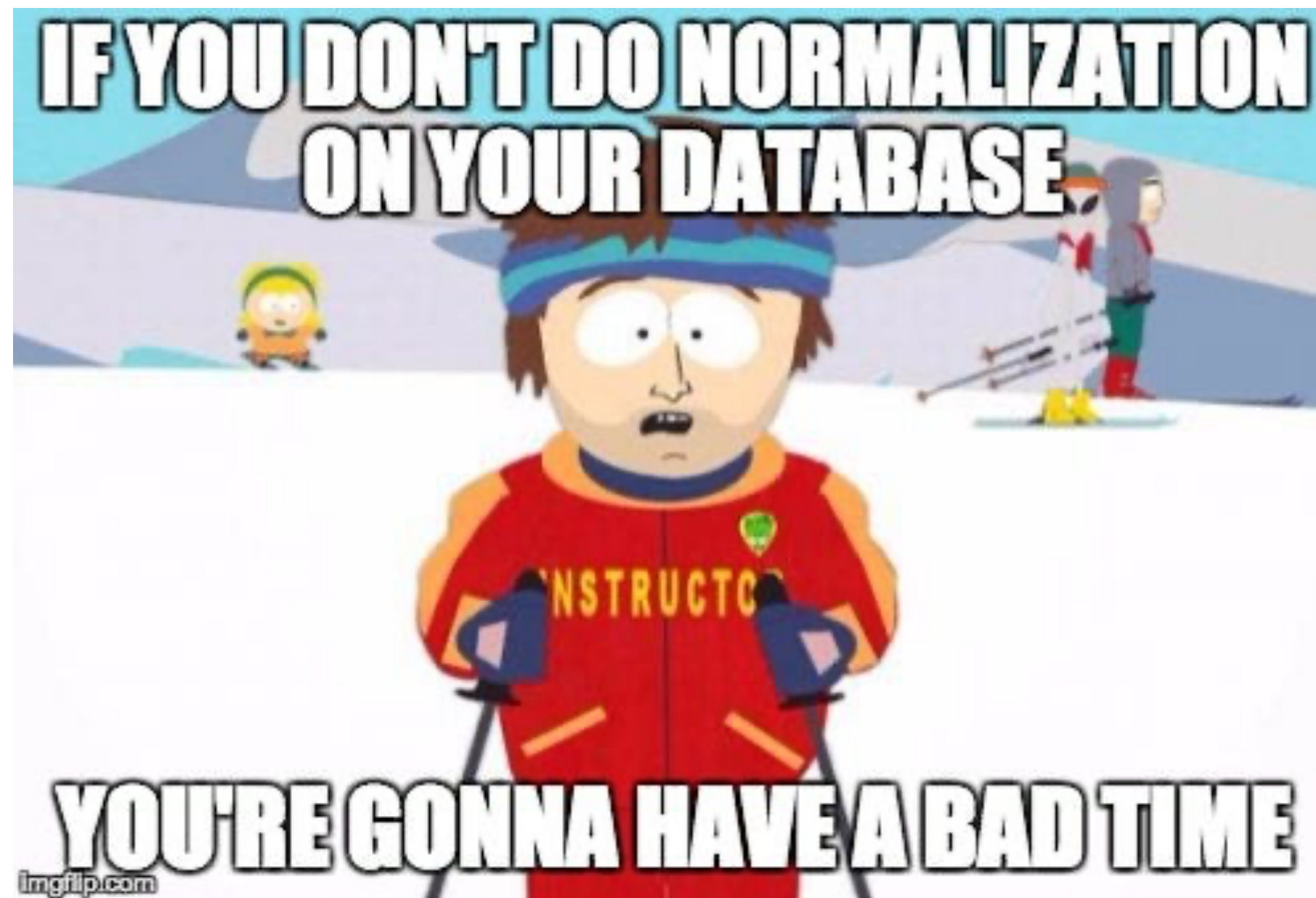
## Most databases in third normal form are fully normalized

There are some exceptions, involving fourth and fifth normal forms

However, these are unusual and outside the scope of what we will cover

# Summing up normalization

Every non-PK column must depend on the PK, the entire PK, and nothing but the PK



# Database design

Figure out what information is fundamental and what is in a supporting role

Bird database: field observations are fundamental, the entire reason for its existence

Each observation is a separate unit and gets its own identifier

This is the primary key of the most fundamental table

Several pieces of information are directly tied to specific observations

E.g., species, time of day, date, location

This information belongs in the observations table

Other pieces of information are independent of individual observations

E.g., conservation status of the species seen, which family it belongs to

This information belongs in other tables

Those form relations and define the primary keys of the supporting tables



## Design your own database, part 1

Design a database to store the results of a set of camera-traps at a field site. Assume there are 12 cameras, each of which may capture from zero to many images per night, that the study lasts for 20 days, and that the study aims to record every case where any kind of mammal sets off the trap. Experts then record what kind of animal was imaged.

The database holds metadata about the images captured by the cameras, and should allow various kinds of queries: by camera, by date, by time, by animal, etc.

## Design your own database, part 2

Sketch a set of tables and columns that could hold this information. Make sure your database is in third normal form. Think about which columns need constraints and whether any foreign key constraints would be helpful.

Experiment by creating tables and loading some sample data. Write code to create the tables and columns for your database. Insert two or more rows of sample data into each table. Design a few test queries to get a feel for how you might retrieve data. Consider making changes and re-create tables as needed. Iterate through changes and testing.





