



Session 2: Ecotoxicological Databases

Theoretical Ecotoxicology

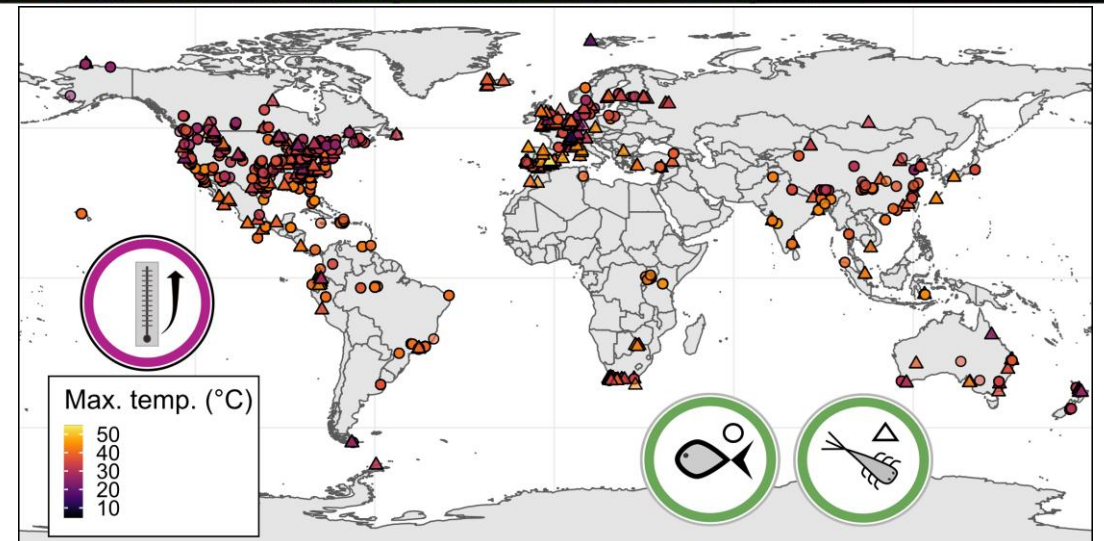
Helena S. Bayat

UNIVERSITÄT
DUISBURG
ESSEN

Offen im Denken

About me

- Helena S. Bayat
- BSc in Environmental Toxicology at UC Davis
 - » Mass spectrometry of aged smoke compounds
- MSc in Environmental Science at University of Copenhagen
 - » Modelling mixture toxicity under stressful food regime
- Currently doing PhD with RESIST project
 - » Multiple stressors in freshwater
 - » Built a database for thermal tolerance



Learning objectives

- Understand what a database is and relevance to ecotoxicology
- Introduce a selection of ecotoxicological databases
- Know where to look for ecotoxicological information

What are databases?

What are databases?

Why do we need them?



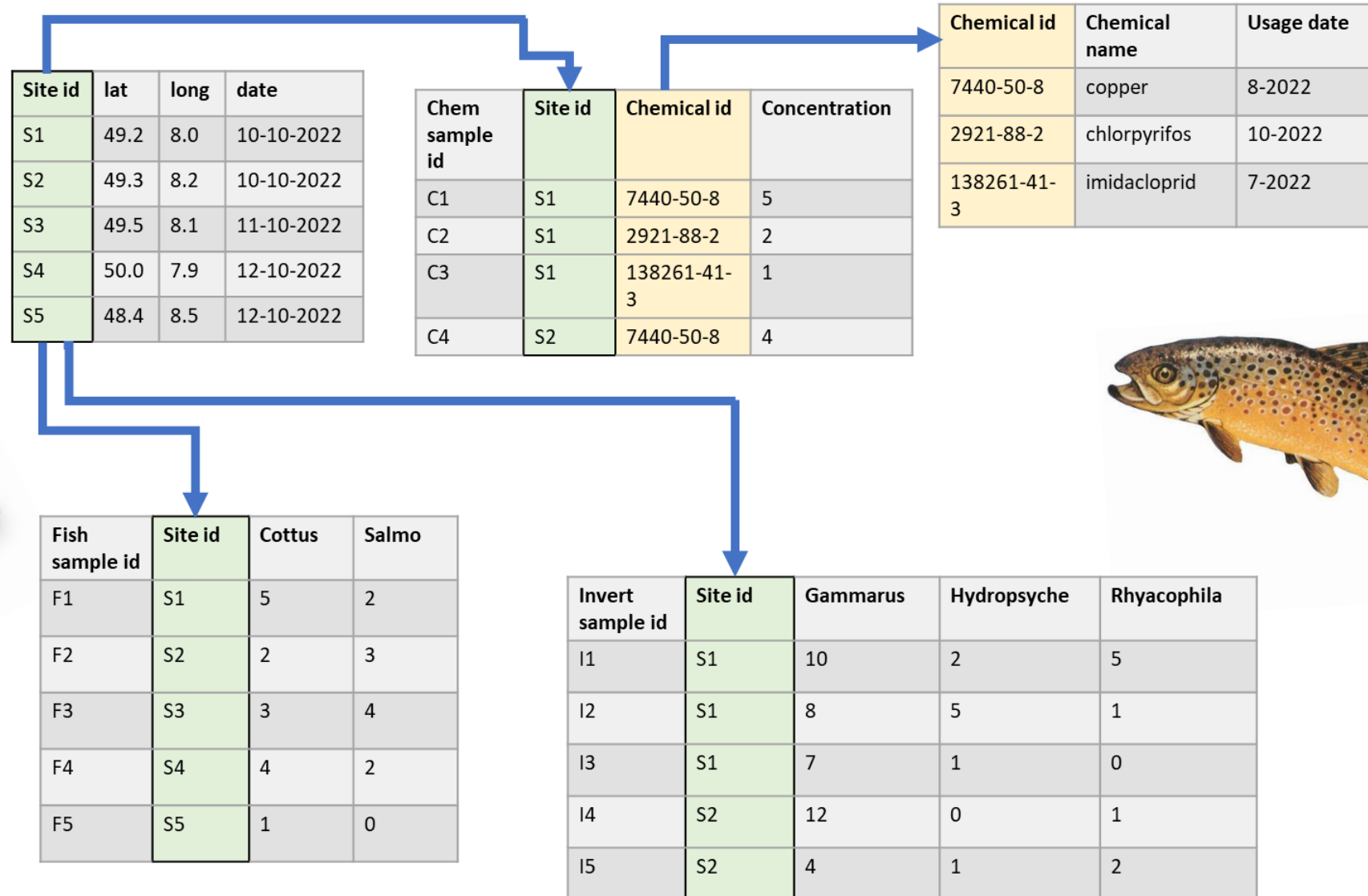
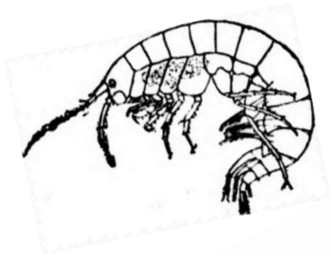
PostgreSQL

Consider:

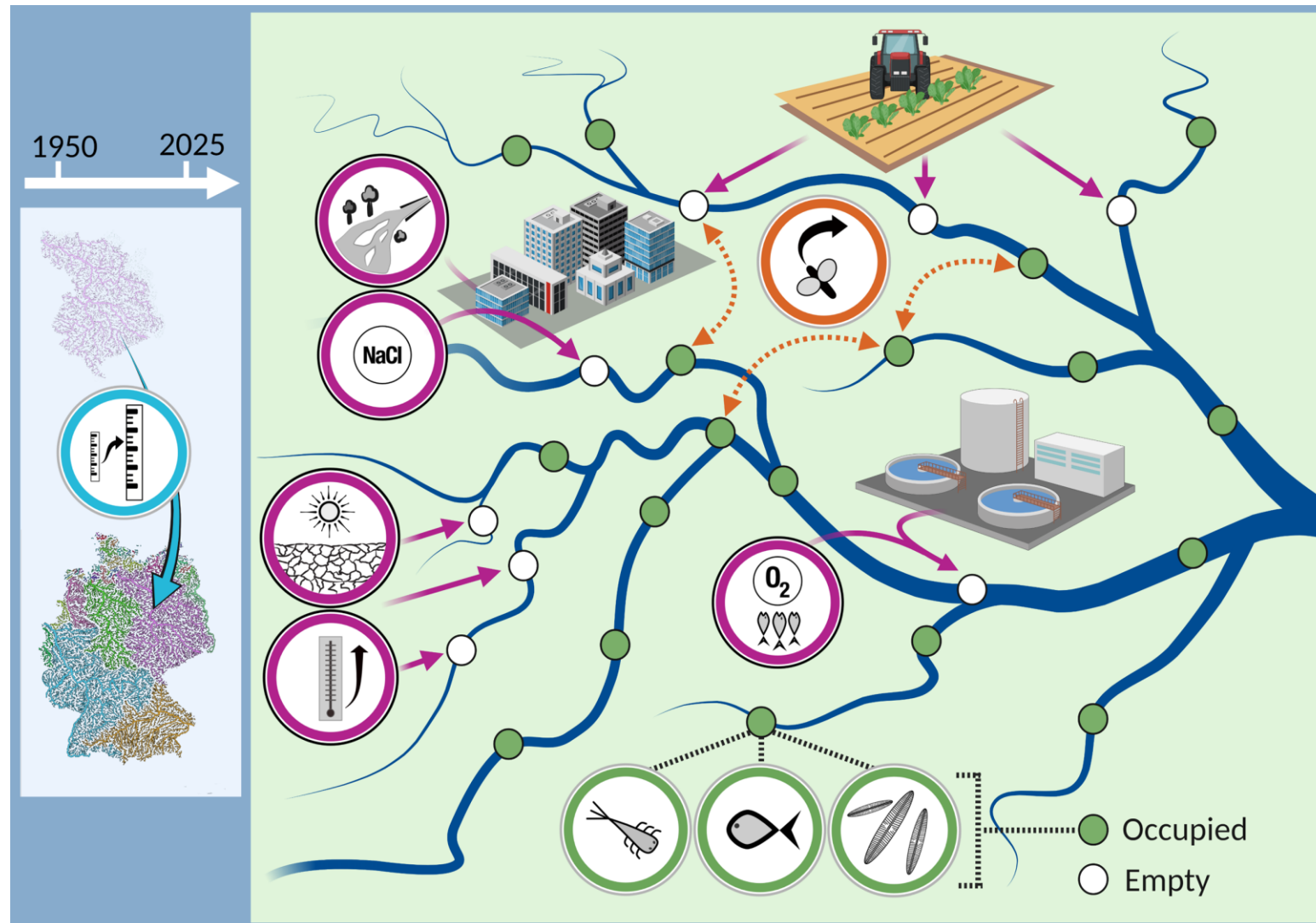
- You have 5 sampling sites
 - » At each you have taken biological samples
 - › Invertebrates
 - › Fish
 - » You also have chemical data for each
 - › Chemical concentration

- » How would you collect and store all this different data?

Consider:



Consider:



Relational databases

- Relational database management systems (RDMS)
- Developed by Edgar F. Codd for IBM in 1970

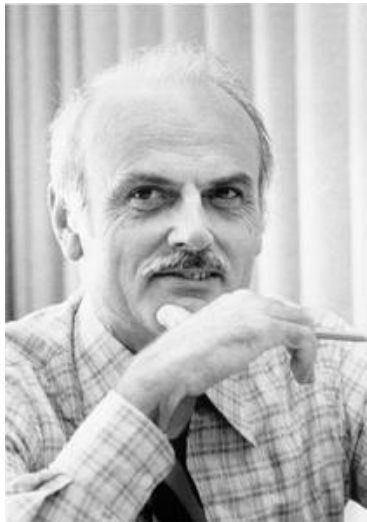


Image from IBM

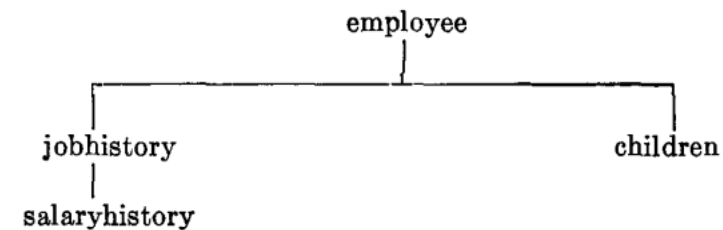
Information Retrieval

P. BAXENDALE, Editor

A Relational Model of Data for Large Shared Data Banks

E. F. CODD
IBM Research Laboratory, San Jose, California

The relational view (or model) of data described in Section 1 appears to be superior in several respects to the graph or network model [3, 4] presently in vogue for non-inferential systems. It provides a means of describing data with its natural structure only—that is, without superimposing any additional structure for machine representation purposes. Accordingly, it provides a basis for a high level data language which will yield maximal independence be-



employee (*man#*, name, birthdate, jobhistory, children)
 jobhistory (*jobdate*, title, salaryhistory)
 salaryhistory (*salarydate*, salary)
 children (*childname*, birthyear)

FIG. 3(a). Unnormalized set

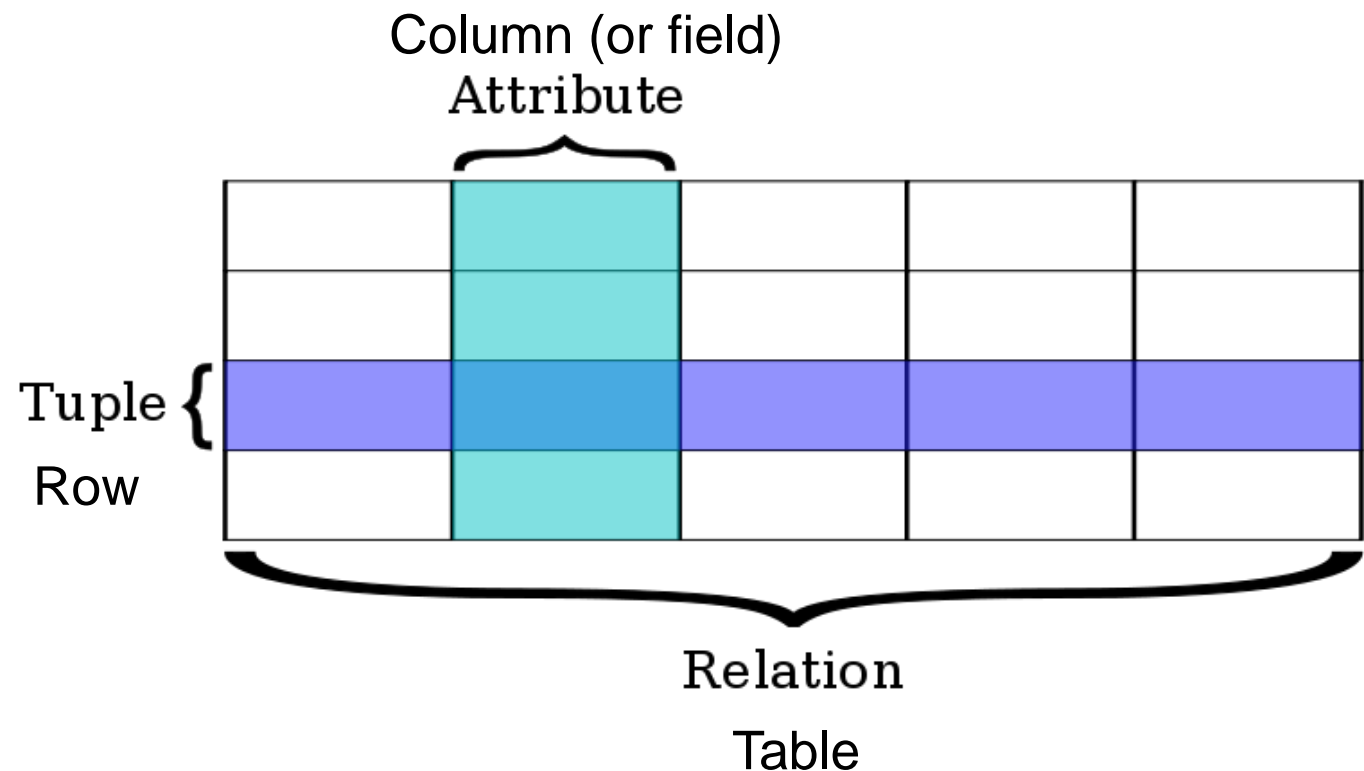
employee' (*man#*, name, birthdate)
 jobhistory' (*man#*, *jobdate*, title)
 salaryhistory' (*man#*, *jobdate*, *salarydate*, salary)
 children' (*man#*, *childname*, birthyear)

FIG. 3(b). Normalized set

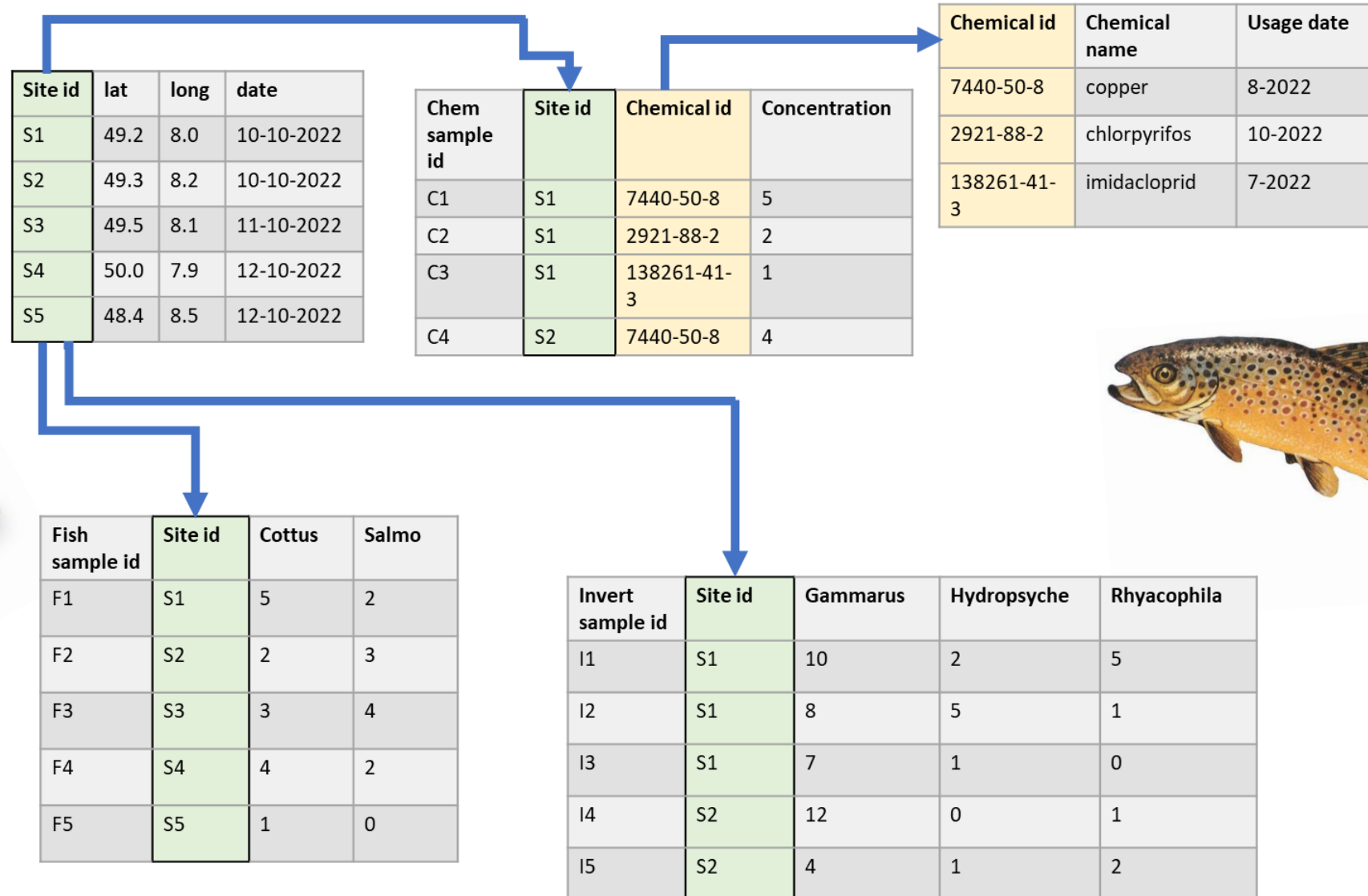
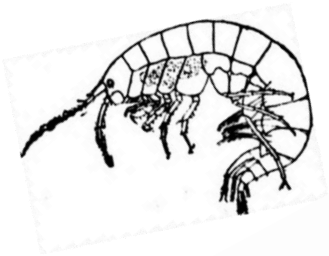
Relational databases

- Data is accessed using Structured Query Language (SQL)
 - » Standardized in 1987, ISO 9075:1987
- Tables connected by keys

```
database.schema.table.column
```



Consider:



Relational databases

SQL software

Free and open source (FOSS)



Proprietary

ORACLE®



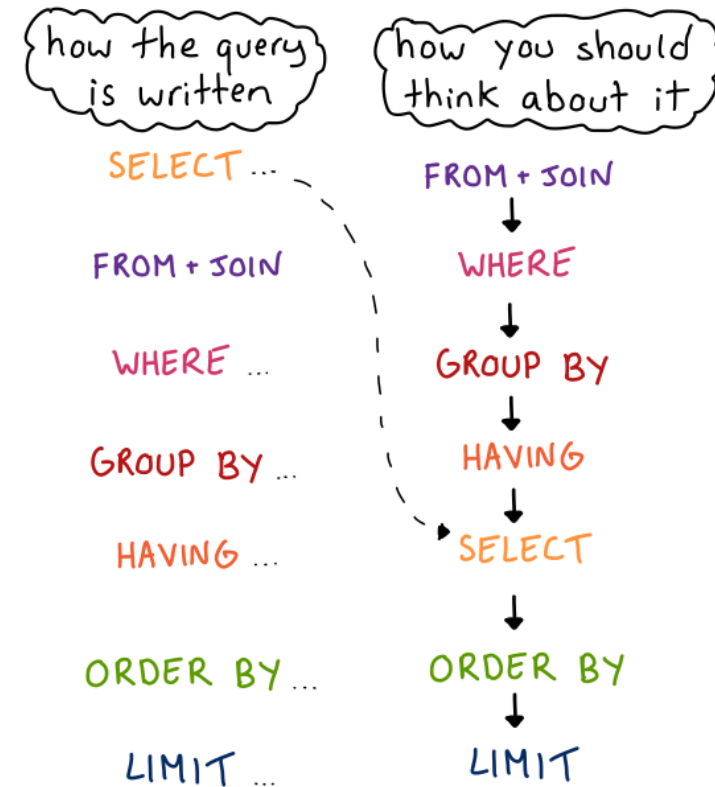
Relational databases

SQL queries

- Query: instructions for retrieving subsets of data

```
SELECT
  city_name,
  country,
  population
FROM cities
WHERE country = 'Austria'
```

The query's steps don't happen in the order they're written:



(In reality query execution is much more complicated than this.
There are a lot of optimizations.)

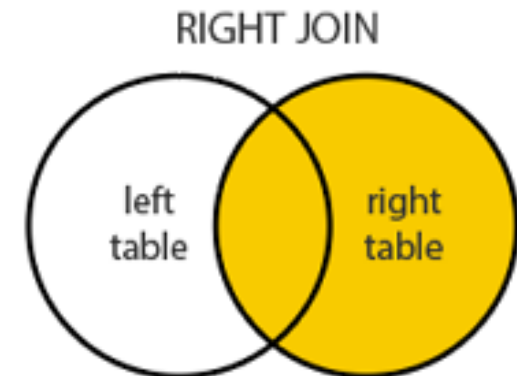
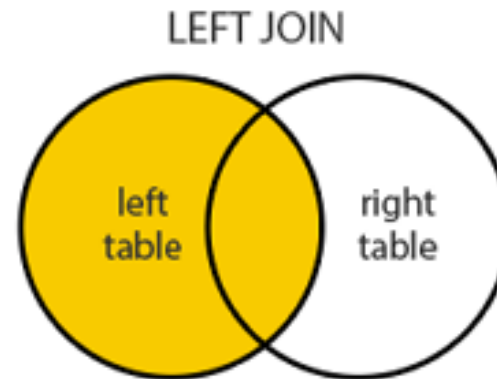
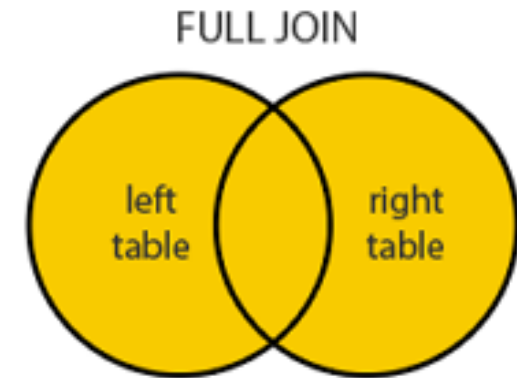
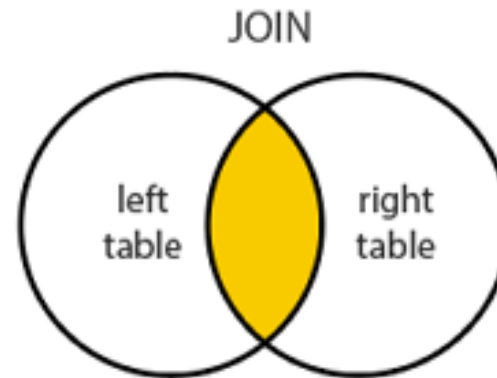
5

SULIA EVANS
@bork

Relational databases

SQL queries

- INNER JOIN
- LEFT JOIN
- RIGHT JOIN
- FULL JOIN



Relational databases

SQL vs R

SQL

- Used with databases
- Data on disk
- Unlimited rows, 250-1600 columns
- Can be slow

R

- Programming language
- Data in memory
- Limited but fast
- Can use packages to connect to databases and send SQL queries

Relational databases

SQL logic in R

- Joins
- Selecting, filtering, grouping, etc
- Part of the tidyverse and data.table

SELECT ...

FROM + JOIN

WHERE ...

GROUP BY ...

HAVING ...

ORDER BY ...

LIMIT ...



Relational databases

SQL logic in R

- Joins
- Selecting, filtering, grouping, etc
- Part of the tidyverse and data.table

For more information:

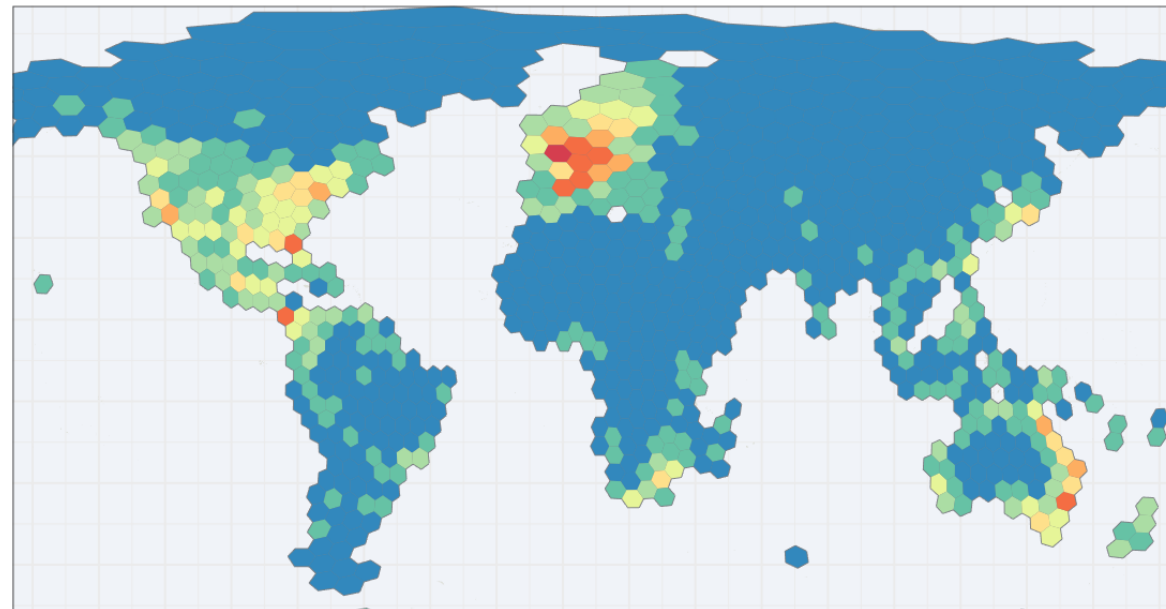
R for Data Science
<https://r4ds.hadley.nz/>

data.table
<https://cran.r-project.org/web/packages/data.table/vignettes/datatable-intro.html>

SELECT ...
FROM + JOIN
WHERE ...
GROUP BY ...
HAVING ...
ORDER BY ...
LIMIT ...



Large databases for ecological data



freshwaterecology.info



EPA ECOTOX Database

- Large!
- Ecotoxicological test results from published studies and grey literature
- Automatically updated
- Hosted by the United States Environmental Protection Agency (USEPA)



Some reminders for working in R

- Working directory
- Make a folder → keep code organized
- Make comments
- Troubleshooting
 - » Check if working directory is set to correct folder
 - » Check spelling
 - » Extra or wrongly placed dots, dashes, etc cause errors in the code
 - » Use Stack Overflow or AI tools when stuck



Time for an exercise