
DataJoint: Bringing databases back into data science

Raphael Guzman
Dimitri Yatsenko, PhD



Intro

- developer
- technologist
- entrepreneur
- oss advocate



<https://github.com/guzman-raphael>



Raphael Guzman



raphael.h.guzman@gmail.com

<https://www.linkedin.com/in/guzman-raphael/>

—

How many of you are using relational databases?

e.g. MySQL, Postgres, SQLite

—

How many don't use relational databases because of **SQL**?

Challenge

Reproducible compute
pipelines in neuroscience
and access to their
produced results



"neuroscientist doing
neuroscience when
neuroscience needs to
be done"



"A friendly database that
helps scientists"

What do databases do?

- **Data structure** reflects the logic of scientific study
- **Data integrity:** entity integrity, completeness, data types, referential integrity, group integrity, transactions
- **Consistency:** everybody sees the same – concurrent access, transactions
- **Precise queries:** get exact slice of data required for analysis
- **Performance:** indexed lookup
- **Support computational dependencies!**

Database revolutions

1951: Magnetic Tape
1955: Magnetic Disk
1961: ISAM
1965: Hierarchical model
1968: IMS
1969: Network Model
1971: IDMS

2003: MarkLogic
2004: MapReduce
2005: Hadoop
2005: Vertica
2007: Dynamo
2008: Cassandra
2008: Hbase
2008: NuoDB
2009: MongoDB
2010: VoltDB
2010: Hana
2011: Riak
2012: Areospike
2014: Splice Machine

1950 - 1972

Pre-Relational

1972 - 2005

Relational

2005 - 2015

The Next Generation

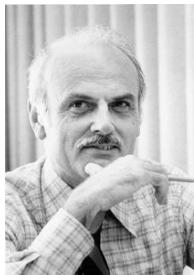
1970: Codd's Paper
1974: System R
1978: Oracle
1980: Commercial Ingres
1981: Informix
1984: DB2
1987: Sybase
1989: Postgres
1989: SQL Server
1995: MySQL

NoSQL
NOSQL
NewSQL

from “Next generation databases: NoSQL and Big Data” by Guy Harrison

Relational Data Model

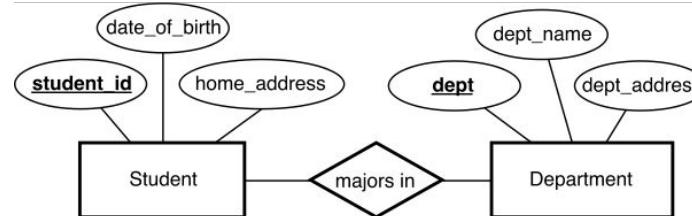
- ❖ Data are in flat simple tables (relations)
- ❖ Domain constraints
- ❖ Unique constraints
- ❖ Referential constraints
- ❖ Declarative queries



Ted F. Codd
IBM, 1969

Entity-Relationship Model

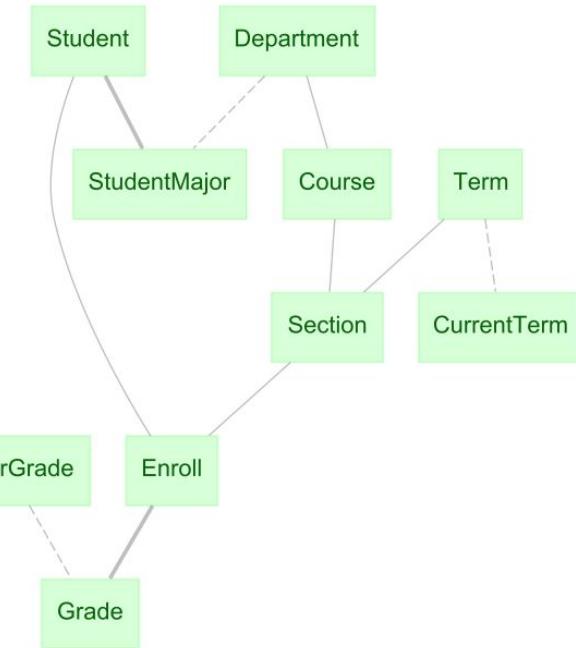
Peter Chen, 1976



- ❖ Conceptual clarification
- ❖ Diagramming notation

DataJoint

<https://arxiv.org/abs/1807.11104>



- ❖ Conceptual refinement
- ❖ Definition language
- ❖ Query language
- ❖ Diagramming notation
- ❖ Integrated computation

What is DataJoint?

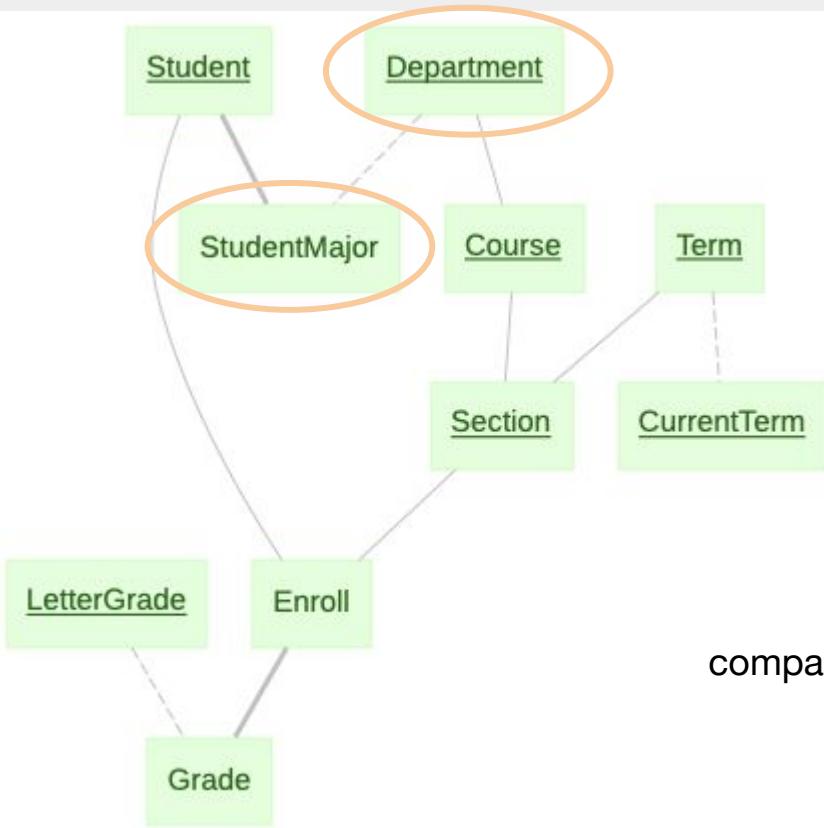
An API to help scientists with:

- Defining computational dependencies
- Querying
- Data lineage



"A researcher unwrangling
the mysteries of the brain"

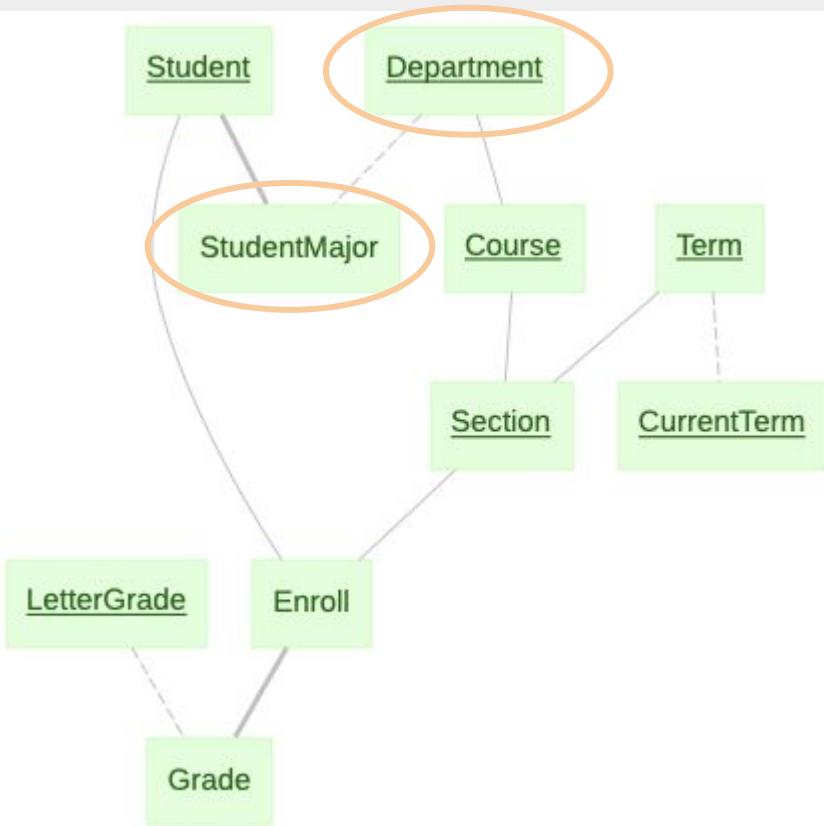
Data Definition Language and diagramming notation



```
CREATE TABLE Department(
    dept char(6) NOT NULL COMMENT "
        abbreviated department name, e.g.
        BIOL",
    dept_name varchar(200) NOT NULL COMMENT "
        full department name",
    dept_address varchar(200) NOT NULL
        COMMENT "mailing address",
    dept_phone varchar(14),
    PRIMARY KEY(dept))

CREATE TABLE StudentMajor(
    student_id int unsigned NOT NULL COMMENT
        "university ID",
    dept char(6) NOT NULL COMMENT "
        abbreviated department name, e.g.
        BIOL",
    declare_date date NOT NULL COMMENT "when
        student declared her major",
    PRIMARY KEY (student_id),
    FOREIGN KEY (student_id) REFERENCES
        Student(student_id),
    FOREIGN KEY (dept) REFERENCES Department(
        dept))
```

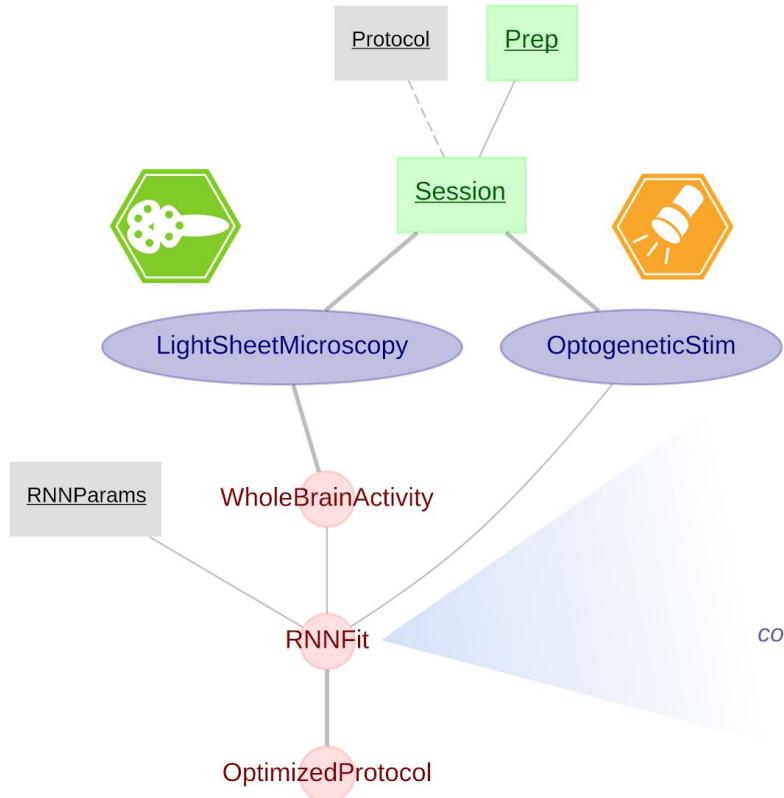
Data Definition Language and diagramming notation



```
: :Department
dept : char(6)      # abbreviated department name, e.g. BIOL
---
dept_name      : varchar(200)    # full department name
dept_address   : varchar(200)    # mailing address
dept_phone     : varchar(14)

: :StudentMajor
-> Student
---
-> Department
declare_date : date    # when student declared her major
```

Computation is part of the data model



database link table name

`@schema
class RNNFit(dj.Computed):
 definition = """
 -> WholeBrainActivity
 -> OptogeneticStim
 -> RNNParams
 ...
 fit : longblob
 cv : longblob # cross-validation
 score : double # fit quality score
 """`

dependency

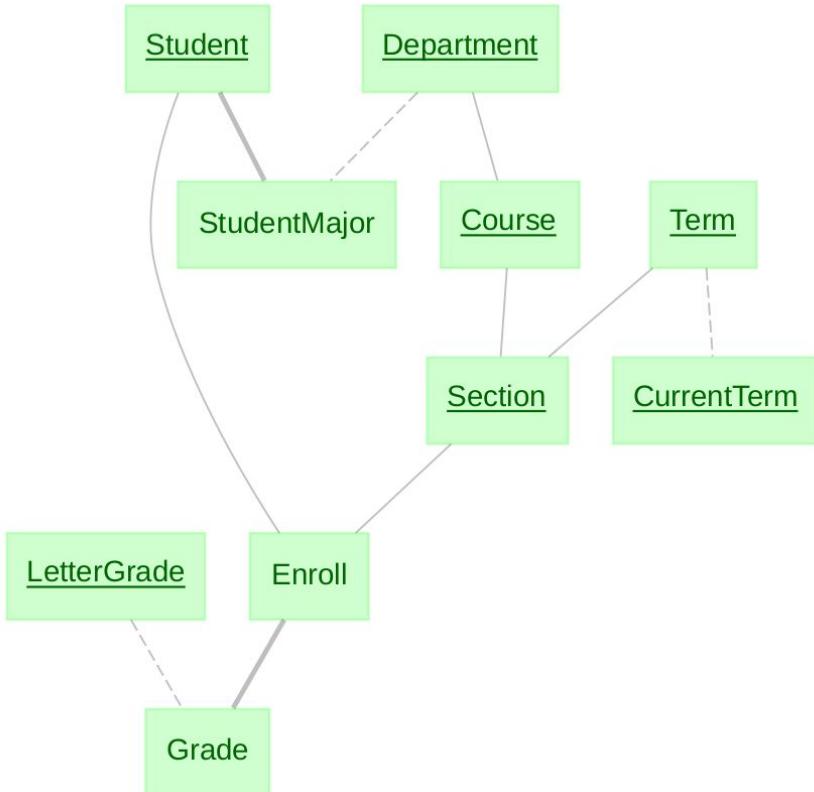
`def make(self, key):
 # fetch data
 params = (RNNParams & key).fetch()
 neural_activity = (WholeBrainActivity & key).fetch()
 opto_stim = (OptogeneticStim & key).fetch()

 # compute
 fit, cv, score = RNN.fit(neural_activity, opto_stim, **params)

 # submit result
 self.insert1(dict(key, fit=fit, cv=cv, score=score))`

computed attributes

Query Language



Query: "Currently enrolled students"

DataJoint:

Student & (**Enroll** * **CurrentTerm**)

SQL:

```
SELECT * FROM Student  
WHERE student_id IN (  
    SELECT student_id FROM Enroll  
    WHERE (term_year, term) IN (  
        SELECT term_year, term  
        FROM CurrentTerm))
```

Live Demo?

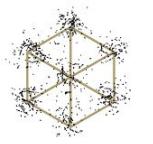


Community of 60+ Research Labs

<https://datajoint.com/docs/projects>



COLUMBIA | ZUCKERMAN INSTITUTE
Mortimer B. Zuckerman Mind Brain Behavior Institute



Kavli Institute
for Systems Neuroscience



UC San Diego
UCLA



Caltech
EPFL



Ψ INDIANA UNIVERSITY



Northwestern
University



W UNIVERSITY of WASHINGTON



Sainsbury Wellcome Centre

Key Contributors

Dimitri Yatsenko

Edgar Y. Walker

Fabian Sinz

Jacob Reimer

Andreas Hoenselaar

Alex Ecker

Philipp Berens

Manolis Froudarakis

Raphael Guzman

Kabilar Gunalan

Thinh Nguyen

Chris Turner

Adib Baji

Jeroen Verswijver

Shan Shen

Any many more...



"superheroes that improve science through code"

Thank you!

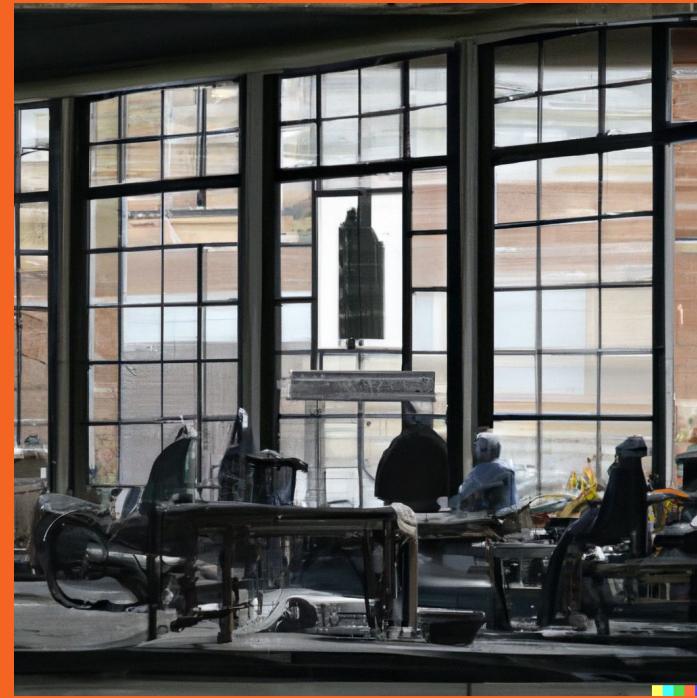


Sponsored by NIH Grant U24 NS116470



Goals

- We seek to elevate DataJoint as an integral part of the SciPy data science stack
- a general tool for scientific relational database programming
- .. with support for computational dependencies
- solid theoretical foundations
- community-driven open-source project
- bolster open-source governance, community feedback, contribution guidelines
- engage NumFocus and other open-source advocates



“fostering a community of respect and acceptance”



Raphael Guzman



Community

- <https://datajoint.com/docs>
- <https://github.com/datajoint/datajoint-tutorials>
- <https://stackoverflow.com/questions/tagged/datajoint>

