

# SQL 101

Ethan Kent

Spoonflower

December 17, 2020

# What is SQL?

- ▶ SQL is a programming language for interacting with relational databases.
- ▶ It stands for *Structured Query Language*.<sup>1</sup>
- ▶ It was invented in 1974, and became a standard of ANSI and ISO in the late-1980s.
- ▶ Standards for SQL are now defined by ISO/IEC, though it appears that no actual SQL versions are compliant with all and only the ISO/IEC standards.

---

<sup>1</sup>Or maybe it doesn't: "It's a common misconception that *SQL* stands for *structured query language*; it stands for S-Q-L and nothing else. Why? Because ISO says so." Chris Fehily, *SQL Database Programming* at loc'n 222 (2015 Kindle ed.)

# What is a Relational Database?

- ▶ A database (imperfectly) implements a *relational model* of data, which is a mathematical model describing certain data structures and relationships.
- ▶ From an practical engineering perspective, a relational database has rows and columns, where each row is an entry in the database; and each column is a bit of data about that entry that matches the type required for that column.
- ▶ A relational database contrasts with NoSQL databases like MongoDB, in which users can insert free-form data that does not comprise rows of data in typed columns.

# What's up with MySQL, PostgreSQL, Microsoft SQL Server, Oracle, IBM DB2, ...

- ▶ The situation is similar to the existence of many programming languages aimed at the same problem domains.
- ▶ Some are open-source, some aren't.
- ▶ Some are more compliant with ISO/IEC standards (e.g. PostgreSQL); some are less so (e.g. SQLite).
- ▶ Some are part of ecosystems, are cloud-based/SaaS, etc. (e.g. Microsoft SQL Server, Amazon RDS, Google Cloud SQL).

# What does “relational” mean?

## Spreadsheet analogy

You’ve probably encountered this problem if you’ve used a spreadsheet to keep track of things. It starts out simple:

Person	Age
Ethan	39
Madigan	39
Belle	9
Milo	7
Clara	3

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

Now we add a bit:

Person	Age	Favorite Dessert
Ethan	39	Oreo Cookies
Madigan	39	Klondike Bar
Belle	9	Mint Chip Ice Cream
Milo	7	Cookie Dough Ice Cream
Clara	3	Chocolate Ice Cream

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

And now it turns out one person can't decide and has *two* favorites.  
We could do this:

Person	Age	Favorite Dessert
Ethan	39	Oreo Cookies, Mint Chip Ice Cream
Madigan	39	Klondike Bar
Belle	9	Mint Chip Ice Cream
Milo	7	Cookie Dough Ice Cream
Clara	3	Chocolate Ice Cream

Or maybe we try this:

Person	Age	Favorite Dessert #1	Favorite Dessert # 2
Ethan	39	Oreo Cookies	Mint Chip Ice Cream
Madigan	39	Klondike Bar	
Belle	9	Mint Chip Ice Cream	
Milo	7	Cookie Dough Ice Cream	
Clara	3	Chocolate Ice Cream	

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

What's wrong with this?

- ▶ The comma-separated list is bad because—
  - ▶ We can't sort or group things together, and
  - ▶ The value has to be manipulated to be used.
- ▶ The use of additional columns is bad because—
  - ▶ We have to keep adding columns based on the data, so we're never sure we're done with the structure of the database; and
  - ▶ We are treating like things differently: columns mean the same thing, but are distinct.



# What does “relational” mean?

Spreadsheet analogy (cont'd)

So we do this:

Person	Age	Favorite Dessert
Ethan	39	Oreo Cookies
Ethan	39	Mint Chip Ice Cream
Madigan	39	Klondike Bar
Belle	9	Mint Chip Ice Cream
Milo	7	Cookie Dough Ice Cream
Clara	3	Chocolate Ice Cream

# What does “relational” mean?

Spreadsheet analogy (cont'd)

And now we want to have more detail about the desserts:

Person	Age	Favorite Dessert	Dessert Type
Ethan	39	Oreos	Cookie
Ethan	39	Mint Chip	Ice Cream
Madigan	39	Klondike Bar	Ice Cream Novelty
Belle	9	Mint Chip	Ice Cream
Milo	7	Cookie Dough	Ice Cream
Clara	3	Chocolate	Ice Cream

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

At this point we might feel uneasy:

- ▶ We're repeating ourselves—
  - ▶ Once we know the person is *Ethan*, we know we have a 39-year-old person.
  - ▶ Once we know the dessert is *Mint Chip*, we know we have an *Ice Cream*.
- ▶ We have have data about different stuff in the same place.

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

Repeating ourselves is more than an inconvenience. Consider this:

Person	Age	Favorite Dessert	Dessert Type
Ethan	39	Oreos	Cookie
Ethan	38	Mint Chip	Ice Cream
Madigan	39	Klondike Bar	Ice Cream Novelty
Belle	9	Mint Chip	Candy Bar
Milo	7	Cookie Dough	Ice Cream
Clara	3	Chocolate	Ice Cream

- ▶ Are these two desserts with the same name (“Mint Chip”)?
- ▶ Are there two *Ethans*?
- ▶ Are there mistakes?
- ▶ If there are mistakes, is Ethan 38 or 39? Is Mint Chip an Ice Cream or a Candy Bar?

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

So let's split things into different tables. And let's use some arbitrary, unique “ID”s, so we can talk about a particular row without having to use some of the data (that could change):

ID	Person	Age
1	Ethan	39
2	Madigan	39
3	Belle	9
4	Milo	7
5	Clara	3

ID	Dessert	Dessert Type
1	Oreos	Cookies
2	Mint Chip	Ice Cream
3	Klondike Bar	Ice Cream Novelty
4	Cookie Dough	Ice cream
5	Chocolate	Ice Cream

But now how can we show the relationship between these? Wait, did I just say *relationship*? As in *relational database*? Holy cow.

# What does “relational” mean?

Spreadsheet analogy (cont'd)

ID	Person	Age
1	Ethan	39
2	Madigan	39
3	Belle	9
4	Milo	7
5	Clara	3

ID	Dessert	Dessert Type
1	Oreos	Cookies
2	Mint Chip	Ice Cream
3	Klondike Bar	Ice Cream Novelty
4	Cookie Dough	Ice cream
5	Chocolate	Ice Cream

Person ID	Dessert ID
1	1
1	2
2	3
3	2
4	4
5	5

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

Did you notice the error that will bite us later? Cookie Dough is an “Ice cream,” not an “Ice Cream.” Let’s fix that:

ID	Person	Age
1	Ethan	39
2	Madigan	39
3	Belle	9
4	Milo	7
5	Clara	3

ID	Desert Type
1	Cookie
2	Ice Cream
3	Ice Cream Novelty

ID	Dessert	Dessert Type ID
1	Oreos	1
2	Mint Chip	2
3	Klondike Bar	3
4	Cookie Dough	2
5	Chocolate	2

Person ID	Dessert ID
1	1
1	2
2	3
3	2
4	4
5	5

# What does “relational” mean?

## Spreadsheet analogy (cont'd)

What does any of this have to do with anything? Well, you've just seen examples of:

- ▶ Relational database design
- ▶ Foreign keys
- ▶ Primary keys
- ▶ Composite primary keys
- ▶ Join tables
- ▶ First and Third Normal Form



# Primary keys

ID	Person	Age
1	Ethan	39
2	Madigan	39
3	Belle	9
4	Milo	7
5	Clara	3

ID	Dessert Type
1	Cookie
2	Ice Cream
3	Ice Cream Novelty

ID	Dessert	Dessert Type ID
1	Oreos	1
2	Mint Chip	2
3	Klondike Bar	3
4	Cookie Dough	2
5	Chocolate	2

Person ID	Dessert ID
1	1
1	2
2	3
3	2
4	4
5	5

# Foreign keys

ID	Dessert	Dessert Type ID
1	Oreos	1
2	Mint Chip	2
3	Klondike Bar	3
4	Cookie Dough	2
5	Chocolate	2

Person ID	Dessert ID
1	1
1	2
2	3
3	2
4	4
5	5

# Composite primary keys

Person ID	Dessert ID
1	1
1	2
2	3
3	2
4	4
5	5

# Join tables

Person ID	Dessert ID
1	1
1	2
2	3
3	2
4	4
5	5

# First Normal Form

Don't do either of these things (columns must be “atomic” and there mustn't be “repeating groups”):

Person	Favorite Dessert
Ethan	Oreos, Mint Chip
Madigan	Klondike Bar
Belle	Mint Chip
Milo	Cookie Dough

Person	Favorite Dessert #1	Favorite Dessert # 2
Ethan	Oreos	Mint Chip
Madigan	Klondike Bar	
Belle	Mint Chip	
Milo	Cookie Dough	

# Third Normal Form

Don't do this:

- ▶ If we know *Ethan*, then we automatically know 39; and
- ▶ If we know *Mint Chip*, then we automatically know *Ice Cream*:<sup>2</sup>

Person	Age	Favorite Dessert	Dessert Type
Ethan	39	Oreos	Cookie
Ethan	39	Mint Chip	Ice Cream
Madigan	39	Klondike Bar	Ice Cream Novelty
Belle	9	Mint Chip	Ice Cream
Milo	7	Cookie Dough	Ice Cream
Clara	3	Chocolate	Ice Cream

---

<sup>2</sup>In First and Second Normal Forms and no “transitive dependencies.”

Second Normal Form has to do with composite primary keys, but it's impossible with the tables we've looked at so far.

# Create some tables

```
CREATE TABLE people (  
  id SERIAL NOT NULL,  
  name TEXT NOT NULL UNIQUE,  
  age INTEGER NOT NULL,  
  PRIMARY KEY (id)  
);
```

```
CREATE TABLE dessert_types (  
  id SERIAL NOT NULL,  
  name TEXT NOT NULL UNIQUE,  
  PRIMARY KEY (id)  
);
```

## Create some more tables

```
CREATE TABLE desserts (  
  id SERIAL NOT NULL,  
  name TEXT NOT NULL UNIQUE,  
  dessert_type_id INTEGER NOT NULL,  
  PRIMARY KEY (id),  
  FOREIGN KEY (dessert_type_id)  
    REFERENCES dessert_types(id)  
);
```

```
CREATE TABLE people_desserts (  
  person_id INTEGER NOT NULL,  
  dessert_id INTEGER NOT NULL,  
  PRIMARY KEY (person_id, dessert_id),  
  FOREIGN KEY (person_id) REFERENCES people(id),  
  FOREIGN KEY (dessert_id) REFERENCES desserts(id)  
);
```



## Insert some data

```
INSERT INTO people (name, age) VALUES  
  ('Ethan ', 39),  
  ('Madigan ', 39),  
  ('Belle ', 9),  
  ('Milo ', 7),  
  ('Clara ', 3);
```

```
INSERT INTO dessert_types (name) VALUES  
  ('Cookie'),  
  ('Ice Cream'),  
  ('Ice Cream Novelty');
```

## Now what?

```
INSERT INTO desserts (dessert_name, dessert_type_id) VALUES  
— ('Oreos', ???),
```

## Queries!

```
INSERT INTO desserts (name, dessert_type_id) VALUES
('Oreos',
 (SELECT id FROM dessert_types
  WHERE name = 'Cookie')),

('Mint Chip',
 (SELECT id FROM dessert_types
  WHERE name = 'Ice Cream')),

('Klondike Bar',
 (SELECT id FROM dessert_types
  WHERE name = 'Ice Cream Novelty')),

('Cookie Dough',
 (SELECT id FROM dessert_types
  WHERE name = 'Ice Cream')),

('Chocolate',
 (SELECT id FROM dessert_types
  WHERE name = 'Ice Cream'));
```

## Not very DRY. How about a CTE?

WITH

```
    cookie_id AS  
    (SELECT id FROM dessert_types  
     WHERE dessert_type = 'Cookie'),
```

```
    ice_cream_id AS  
    (SELECT id FROM dessert_types  
     WHERE dessert_type = 'Ice Cream'),
```

```
    ice_cream_novelty_id AS  
    (SELECT id FROM dessert_types  
     WHERE dessert_type = 'Ice Cream Novelty')
```

```
INSERT INTO desserts (dessert_name, dessert_type_id) VALUES  
    ('Oreos', (SELECT id FROM cookie_id)),  
    ('Mint Chip', (SELECT id FROM ice_cream_id)),  
    ('Klondike Bar', (SELECT id FROM ice_cream_novelty_id)),  
    ('Cookie Dough', (SELECT id FROM ice_cream_id)),  
    ('Chocolate', (SELECT id FROM ice_cream_id));
```

Not very DRY. How about a CTE?

Not with MySQL though.