Today's Lecture

1. How to communicate DBs with
   ○ x-teams, customers, stakeholders?

2. What are good designs?

# Database Design

- **Database design: Why?**
  - Agree on schema for use cases now (and later)

- **Consider issues such as:**
  - What entities to model
  - How entities are related
  - What constraints exist in the domain
  - How to achieve <u>good</u> designs

- **Several formalisms exist**
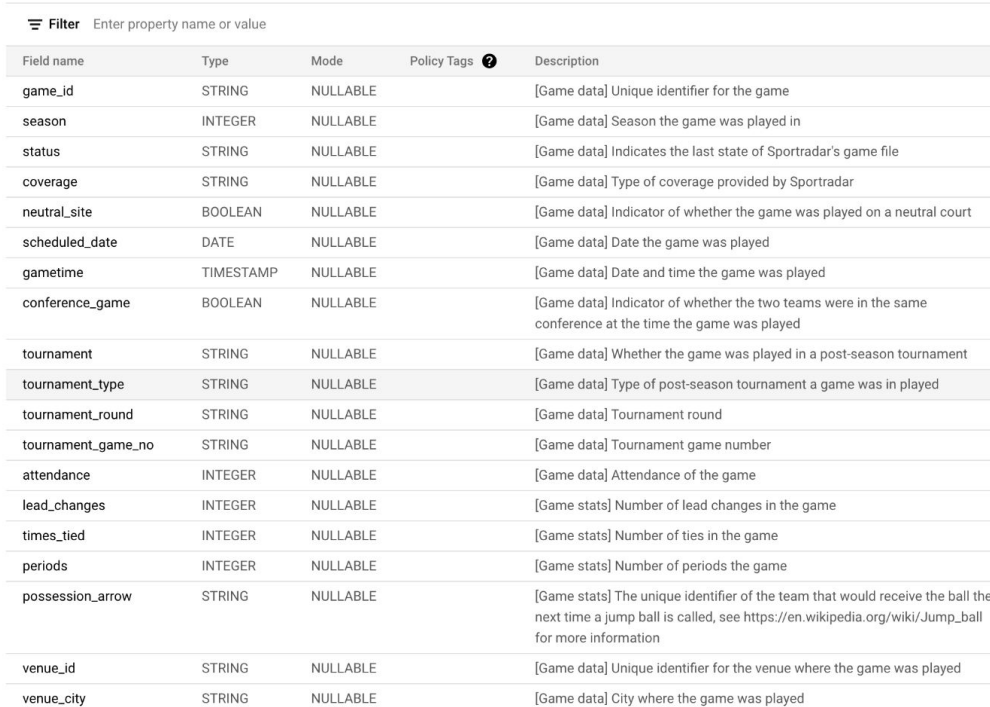  - We discuss some flavors (ER diagrams, DAG diagrams)

# Example 1: NCAA Basketball -- schema for 1 table in BigQuery

| Field name | Type | Mode | Policy Tags ❓ | Description |
|---|---|---|---|---|
| game_id | STRING | NULLABLE | | [Game data] Unique identifier for the game |
| season | INTEGER | NULLABLE | | [Game data] Season the game was played in |
| status | STRING | NULLABLE | | [Game data] Indicates the last state of Sportradar's game file |
| coverage | STRING | NULLABLE | | [Game data] Type of coverage provided by Sportradar |
| neutral_site | BOOLEAN | NULLABLE | | [Game data] Indicator of whether the game was played on a neutral court |
| scheduled_date | DATE | NULLABLE | | [Game data] Date the game was played |
| gametime | TIMESTAMP | NULLABLE | | [Game data] Date and time the game was played |
| conference_game | BOOLEAN | NULLABLE | | [Game data] Indicator of whether the two teams were in the same conference at the time the game was played |
| tournament | STRING | NULLABLE | | [Game data] Whether the game was played in a post-season tournament |
| tournament_type | STRING | NULLABLE | | [Game data] Type of post-season tournament a game was in played |
| tournament_round | STRING | NULLABLE | | [Game data] Tournament round |
| tournament_game_no | STRING | NULLABLE | | [Game data] Tournament game number |
| attendance | INTEGER | NULLABLE | | [Game data] Attendance of the game |
| lead_changes | INTEGER | NULLABLE | | [Game stats] Number of lead changes in the game |
| times_tied | INTEGER | NULLABLE | | [Game stats] Number of ties in the game |
| periods | INTEGER | NULLABLE | | [Game stats] Number of periods the game |
| possession_arrow | STRING | NULLABLE | | [Game stats] The unique identifier of the team that would receive the ball the next time a jump ball is called, see https://en.wikipedia.org/wiki/Jump_ball for more information |
| venue_id | STRING | NULLABLE | | [Game data] Unique identifier for the venue where the game was played |
| venue_city | STRING | NULLABLE | | [Game data] City where the game was played |

Filter ▸ Enter property name or value

Tree navigation:
- ▸ nasa_wildfire
- ▾ ncaa_basketball
  - mascots
  - mbb_games_sr
  - mbb_historical_teams_games
  - mbb_historical_teams_seasons
  - mbb_historical_tournament_games
  - mbb_pbp_sr
  - mbb_pbp_sr-2021-11-14T23_54_22
  - mbb_players_games_sr
  - mbb_teams
  - mbb_teams_games_sr
  - team_colors
- ▸ new_york
- ▸ new_york_311
- ▸ new_york_citibike
- ▸ new_york_mv_collisions
- ▸ new_york_subway
- ▸ new_york_taxi_trips
- ▸ new_york_trees
- ▸ nhtsa_traffic_fatalities
- ▸ nih_gudid
- ▸ nih_sequence_read

1. How did find relationships between columns?
2. How about 10x-100x tables? Columns?

# Example2: Shopify's *simplified* ERD (Entity Relation Diagram)

# Example3: Complex data flows



DBT's DAG for Data Flow



Note: 115 Stages in pipeline !!!

Spark's DAG for Data Flow

# Problems

1. How to connect schemas across

    a. 10s-1000s of Tables, Columns, Relationships, and Flows?

2. How teams collaborate on Big Schemas?
    a. Different subsets useful to App team, Data Analysts, Data Eng...

3. When schemas change?

# Intuition: Cooking Prep

**Source (raw)**

lettuce

onions

cucumber

tomatoes

water

vinegar

Italian spices

**Stage 1**

Chopped lettuce

Sliced onions

Sliced cucumber

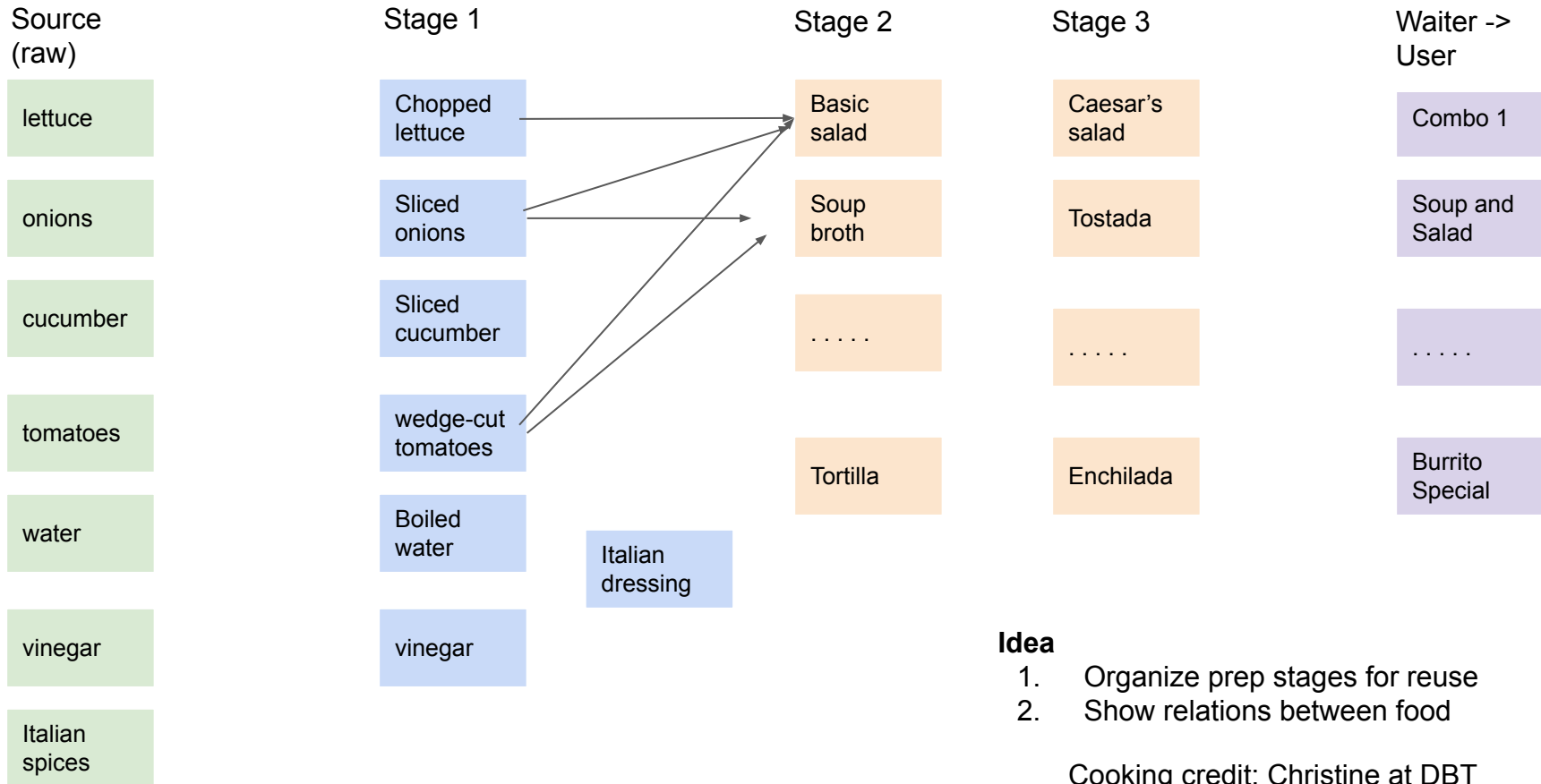wedge-cut tomatoes

Boiled water

Italian dressing

vinegar

**Stage 2**

Basic salad

Soup broth

. . . . .

Tortilla

**Stage 3**

Caesar's salad

Tostada

. . . . .

Enchilada

**Waiter -> User**

Combo 1

Soup and Salad

. . . . .

Burrito Special

**Idea**
1. Organize prep stages for reuse
2. Show relations between food

Cooking credit: Christine at DBT

# Problems



Personas

Schema Graph
(Guha's talk on
datacommons.org)

Big Schemas

# Big Schema

1. Example1: Amazon Product orders


2. Example2: NCAA Basketball schema
    i. "Cooked" version

# For Project3

i. Use tool to convey below. Or something equivalent in text/figures.

ii. Important: **Start from** https://bigschema.io to create new diagram. Don't change the example links.

- **Analysis of your dataset (10%)**
  - +
    - Students show that they are using a meaningful dataset in terms of size and complexity. The overall dataset should be at least 250 MB.
    - Students clearly describe the information captured in the table.
    - It is clear that students understand the structure of their datasets such as data sizes and high-level relationships between tables.
    - Students list the keys and foreign keys between tables that will be used for exploration or describe connections between tables in some other way.
  - -
    - Students use a very simple or very small dataset (e.g. only one table with few columns or a dataset with very few tuples overall).
    - Little to no effort in explaining the dataset. It is not clear to the grader that the student(s) behind the project understands the structure of the data they are working with.