

# Data

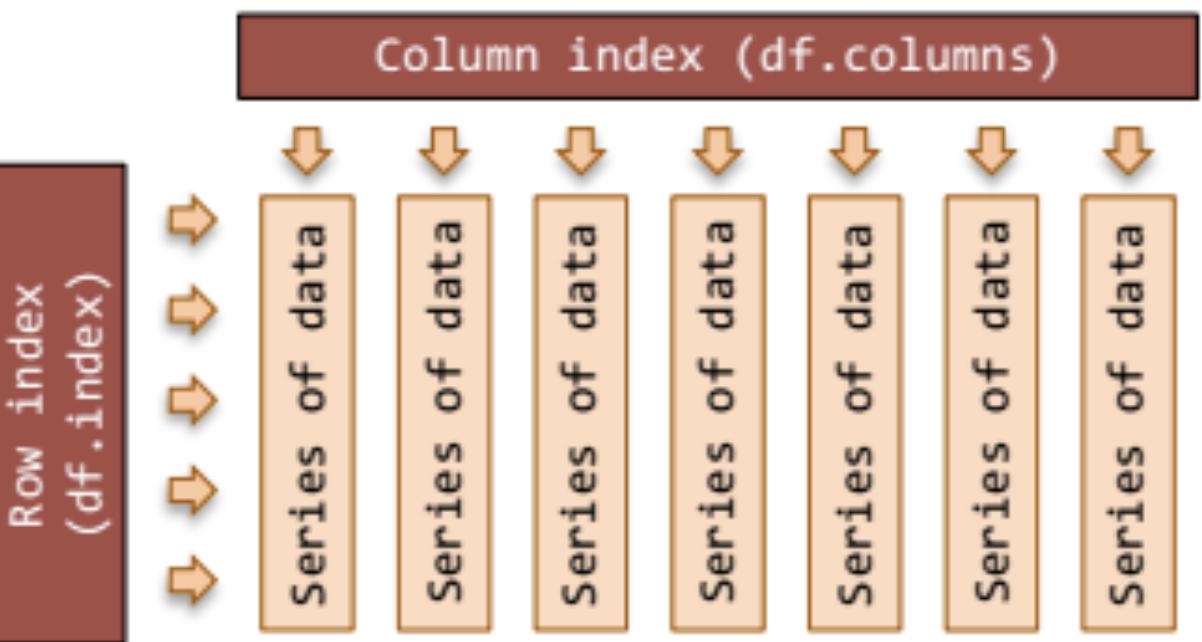
# Databases

# Types in Data

# Tables

Data comes in many ways, but most often it comes in tables. These might be tables in relational databases, Pandas dataframes, csv files on disk, and excel spreadsheets.

**DataFrame object:** a two-dimensional table of data with column and row indexes. The columns are made up of pandas Series objects.



item index	quality	type	price
0	High	Toy	20
1	High	Book	5
2	Medium	Craft	12
3	Low	Craft	10

**Items**

customer index	item index	quantity	total
0	0	5	100
1	0	3	60
2	2	4	48
3	3	10	100

**Customers**

# Types

The diagram illustrates two data tables: **Items** and **Customers**. The **Items** table has columns: **item index**, **quality**, **type**, and **price**. The **Customers** table has columns: **customer index**, **item index**, **quantity**, and **total**. Blue arrows point from the text labels on the left to specific cells in the tables:

- An arrow labeled "Ordinal" points to the **quality** column in the **Items** table.
- An arrow labeled "Nominal/Categorical" points to the **type** column in the **Items** table.
- An arrow labeled "Quantitative" points to the **price** column in the **Items** table.
- An arrow labeled "Foreign Key" points to the **customer index** column in the **Customers** table, which contains values 0, 1, 2, and 3. An arrow also points from the **customer index** column in the **Customers** table to the **item index** column in the **Items** table, labeled "Foreign Key".

**Items**

item index	quality	type	price
0	High	Toy	20
1	High	Book	5
2	Medium	Craft	12
3	Low	Craft	10

**Customers**

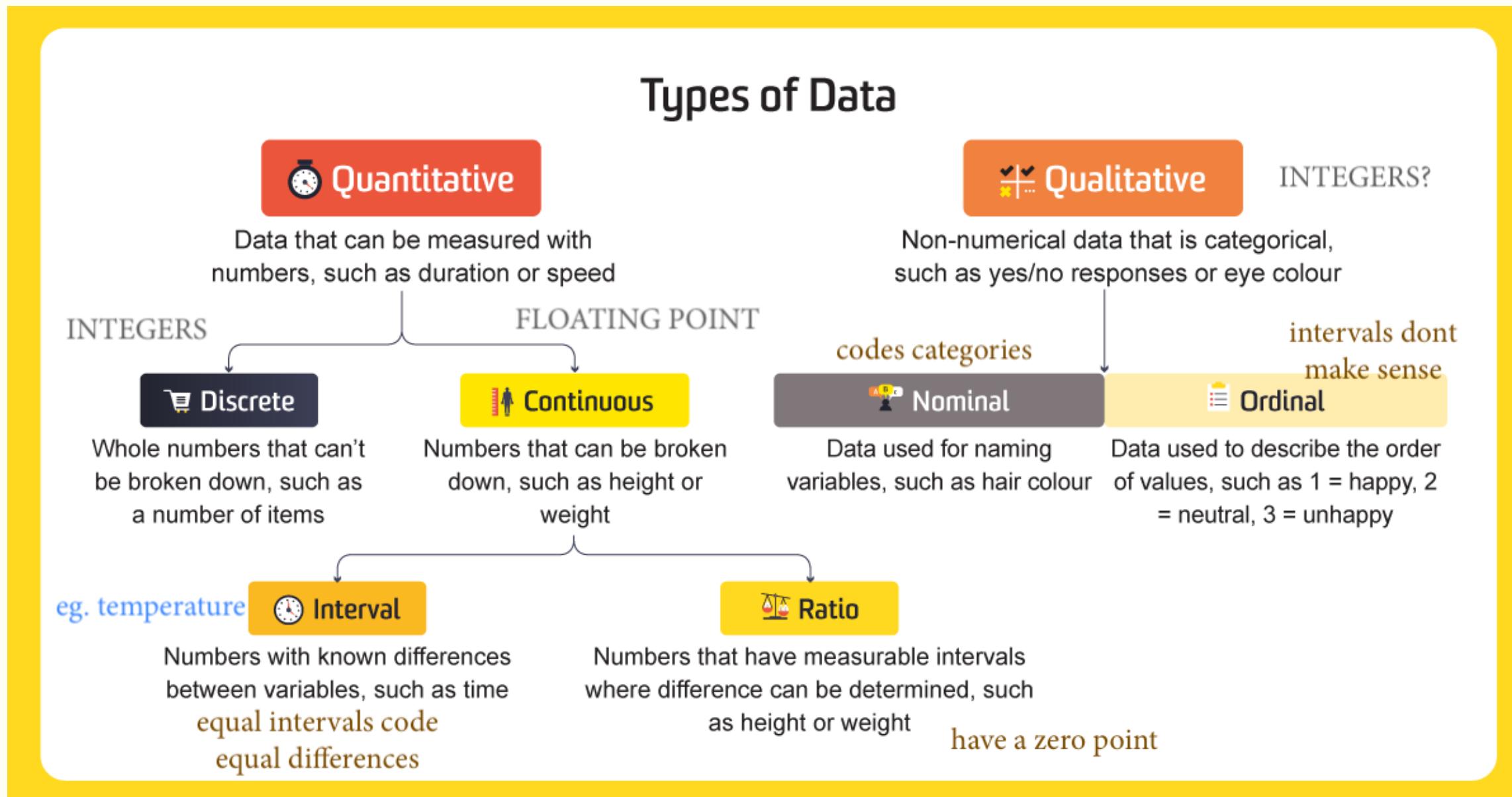
customer index	item index	quantity	total
0	0	5	100
1	0	3	60
2	2	4	48
3	3	10	100

Data has "type". It might be numbers, words, characters, true/false's. And it may be missing, which is a None, or a "NA" or a NULL.

Sometimes data refers to other tables, using a "Foreign Key", such as a reference to another sheet in a spreadsheet, or an id referring to a column in another database table.

Often you must **encode** the data in some form to be useful. For example, you might use integers to encode the ordinal column "quality" here.

# Types: more detail



# Think about the types here!

Table: contributors

New Record    Delete Record

	<a href="#">id</a>	<a href="#">last_name</a>	<a href="#">first_name</a>	<a href="#">middle_name</a>	<a href="#">street_1</a>	<a href="#">street_2</a>	<a href="#">city</a>	<a href="#">state</a>	<a href="#">zip</a>	<a href="#">amount</a>	<a href="#">date</a>	<a href="#">candidate_id</a>
	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter
1	1	Agee	Steven	NULL	549 Laurel ...	NULL	Floyd	VA	24091	500	2007-06-30	16
2	5	Akin	Charles	NULL	10187 Suga...	NULL	Bentonville	AR	72712	100	2007-06-16	16
3	6	Akin	Mike	NULL	181 Baywo...	NULL	Monticello	AR	71655	1500	2007-05-18	16
4	7	Akin	Rebecca	NULL	181 Baywo...	NULL	Monticello	AR	71655	500	2007-05-18	16
5	8	Aldridge	Brittni	NULL	808 Capitol...	NULL	Washington	DC	20024	250	2007-06-06	16
6	9	Allen	John D.	NULL	1052 Cann...	NULL	North Augu...	SC	29860	1000	2007-06-11	16
7	10	Allen	John D.	NULL	1052 Cann...	NULL	North Augu...	SC	29860	1300	2007-06-29	16
8	11	Allison	John W.	NULL	P.O. Box 10...	NULL	Conway	AR	72033	1000	2007-05-18	16
9	12	Allison	Rebecca	NULL	3206 Sum...	NULL	Little Rock	AR	72227	1000	2007-04-25	16

# Encoding types (example SQL)

Data type	Description
CHARACTER(n)	Character string. Fixed-length n
VARCHAR(n) or CHARACTER VARYING(n)	Character string. Variable length. Maximum length n
BINARY(n)	Binary string. Fixed-length n
BOOLEAN	Stores TRUE or FALSE values
VARBINARY(n) or BINARY VARYING(n)	Binary string. Variable length. Maximum length n
INTEGER(p)	Integer numerical (no decimal). Precision p
FLOAT(p)	Approximate numerical, mantissa precision p.
DATE	Stores year, month, and day values
TIME	Stores hour, minute, and second values
TIMESTAMP	Stores year, month, day, hour, minute, and second values
INTERVAL	Representing a period of time, depending on the type of interval

# Representation of Data

	<b>id</b>	<b>first_name</b>	<b>last_name</b>	<b>middle_name</b>	<b>party</b>
	Filter	Filter	Filter	Filter	Filter
1	16	Mike	Huckabee		R
2	20	Barack	Obama		D
3	22	Rudolph	Giuliani		R
4	24	Mike	Gravel		D
5	26	John	Edwards		D
6	29	Bill	Richardson		D
7	30	Duncan	Hunter		R
8	31	Dennis	Kucinich		D
9	32	Ron	Paul		R

Some data is tabular

# And some data is an array of numbers

5
7
4 5
1 2
- 6
3
2 2
1
6
3
- 9

- 9	4	2	5	7
3	0	1 2	8	6 1
1	2 3	- 6	4 5	2
2 2	3	- 1	7 2	6

Images are arrays of numbers.

So are videos.

So is text.

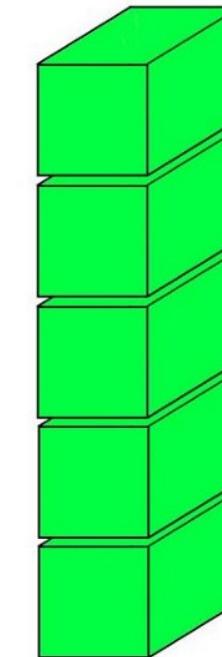
# From One Image



- 9	4	2	5	7	
3	0	1 2	8	6 1	
1	2 3	- 6	4 5	2	
2 2	3	- 1	7 2	6	

(height, width, color\_depth)

# To Many Images



(sample\_size, height, width, color\_depth)

**1 D TENSOR /**  
**VECTOR**

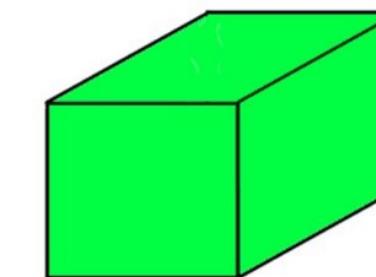
5
7
4 5
1 2
- 6
3
2 2
1
6
3
- 9

**2 D TENSOR /**  
**MATRIX**

- 9	4	2	5	7
3	0	1 2	8	6 1
1	2 3	- 6	4 5	2
2 2	3	- 1	7 2	6

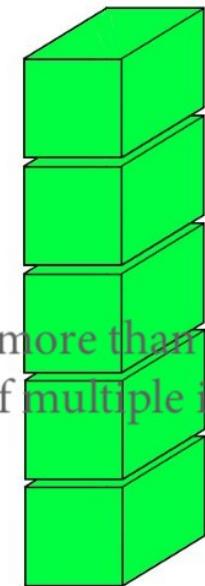
One channel of an image  
a document

**3 D TENSOR /**  
**CUBE**



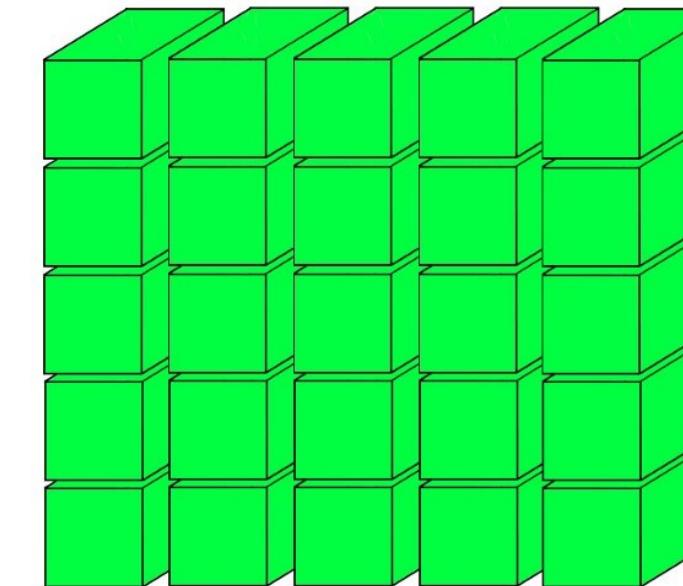
- 9	4	2	5	7
3	0	1 2	8	6 1
1	2 3	- 6	4 5	2
2 2	3	- 1	7 2	6

An image with multiple channels.  
Many 1-channel images  
multiple documents



multiple images with more than one channel  
one video consisting of multiple image frames

**4 D TENSOR**  
**VECTOR OF CUBES**



**5 D TENSOR**  
**MATRIX OF CUBES**

multiple videos

# Text Data

We can store text data in a 3D tensor too. Let's take a look at tweets.

Tweets are max 280 characters. Twitter uses the UTF-8 standard, which allows for millions of types of characters, but let's focus on the first 128, which are basic ASCII.

Then one tweet is a 2D array of shape (280,128). What is in there? What if a tweet is shorter?

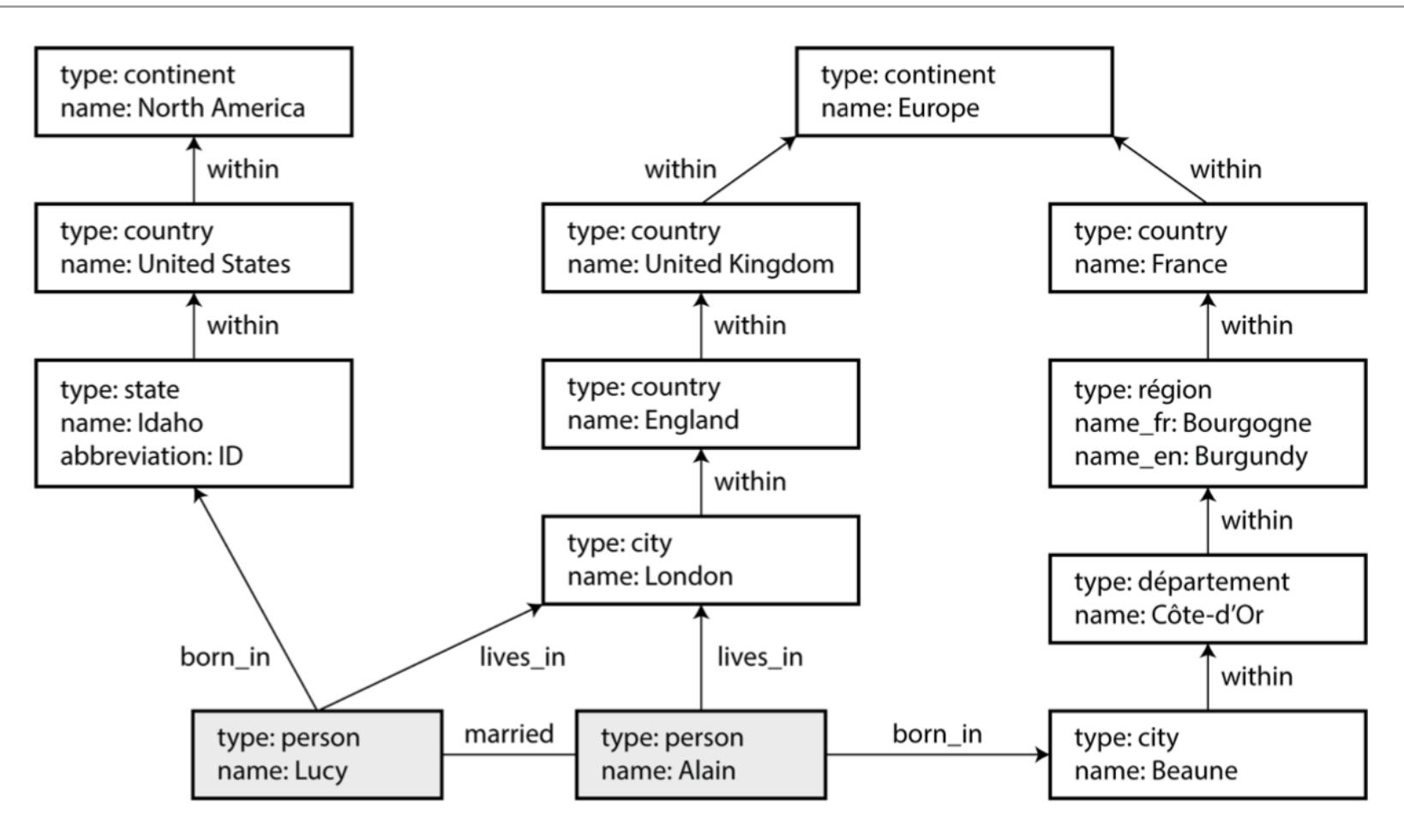
For multiple tweets we go to a 3D array/tensor of shape:

(number\_of\_tweets\_captured, 280, 128)

(As an array this is a wasteful encoding, but there are better formats)

# Graph Data

- Has vertices and edges
- Example: Social connectivity graph, flight paths between airports with edge length being distance (for fuel optimization), computer networks, code graphs, unstructured data
- can be converted into tables or dealt with as is



# Databases and Indexes

<http://www.linkedin.com/in/williamhgates>



#### Summary

Co-chair of the Bill & Melinda Gates Foundation.  
Chairman, Microsoft Corporation. Voracious  
reader. Avid traveler. Active blogger.

#### Experience

Co-chair • Bill & Melinda Gates Foundation  
*2000 – Present*

Co-founder, Chairman • Microsoft  
*1975 – Present*

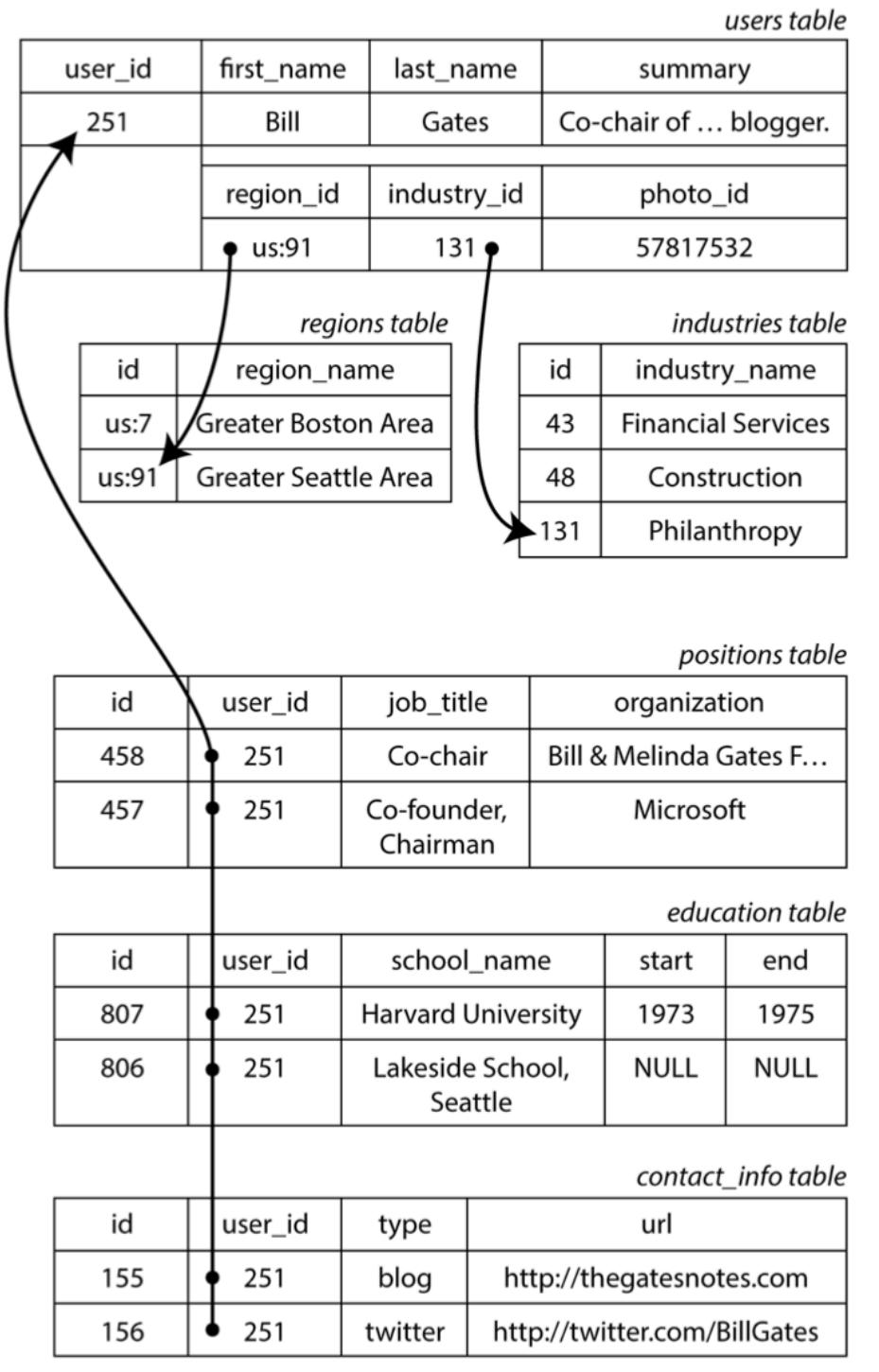
#### Education

Harvard University  
*1973 – 1975*

Lakeside School, Seattle

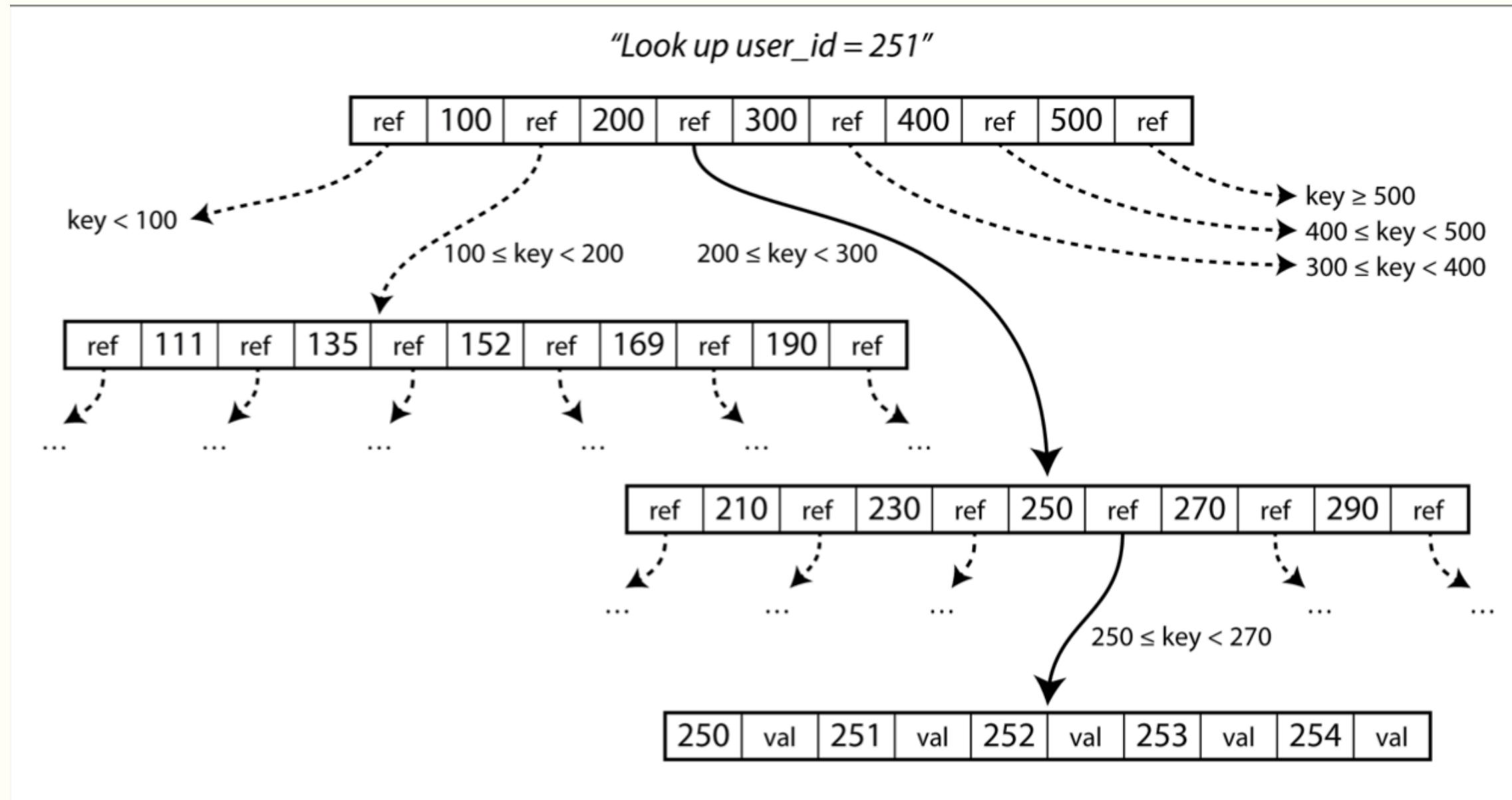
#### Contact Info

Blog: [thegatesnotes.com](http://thegatesnotes.com)  
Twitter: @BillGates



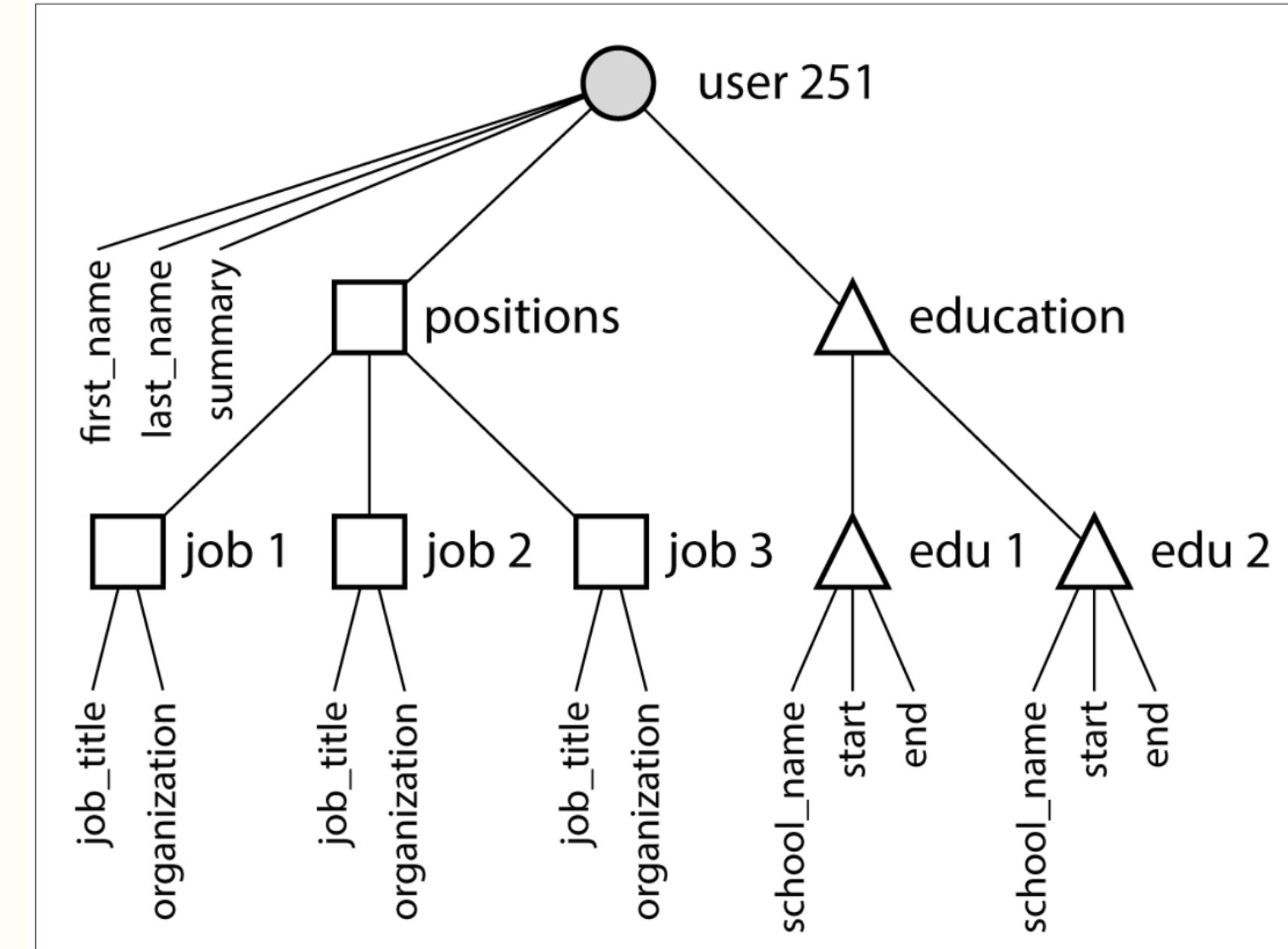
# Relational Database

# How to get stuff?

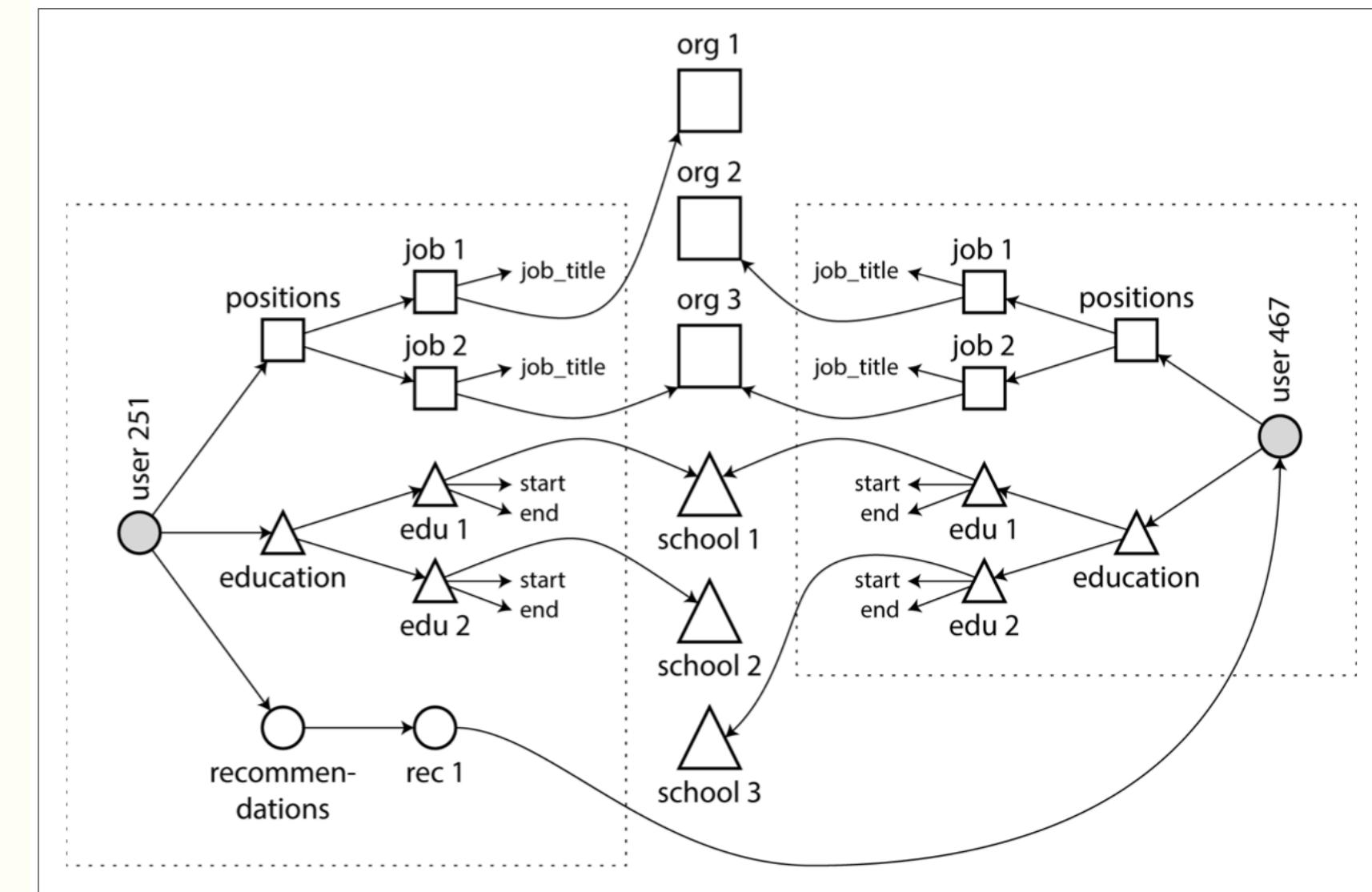


# Document Database

```
{  
  "user_id": 251,  
  "first_name": "Bill",  
  "last_name": "Gates",  
  "summary": "Co-chair of the Bill & Melinda Gates... Active blogger.",  
  "region_id": "us:91",  
  "industry_id": 131,  
  "photo_url": "/p/7/000/253/05b/308dd6e.jpg",  
  "positions": [  
    {"job_title": "Co-chair", "organization": "Bill & Melinda Gates Foundation"},  
    {"job_title": "Co-founder, Chairman", "organization": "Microsoft"}  
  ],  
  "education": [  
    {"school_name": "Harvard University", "start": 1973, "end": 1975},  
    {"school_name": "Lakeside School, Seattle", "start": null, "end": null}  
  ],  
}
```



# Graph or Relational?



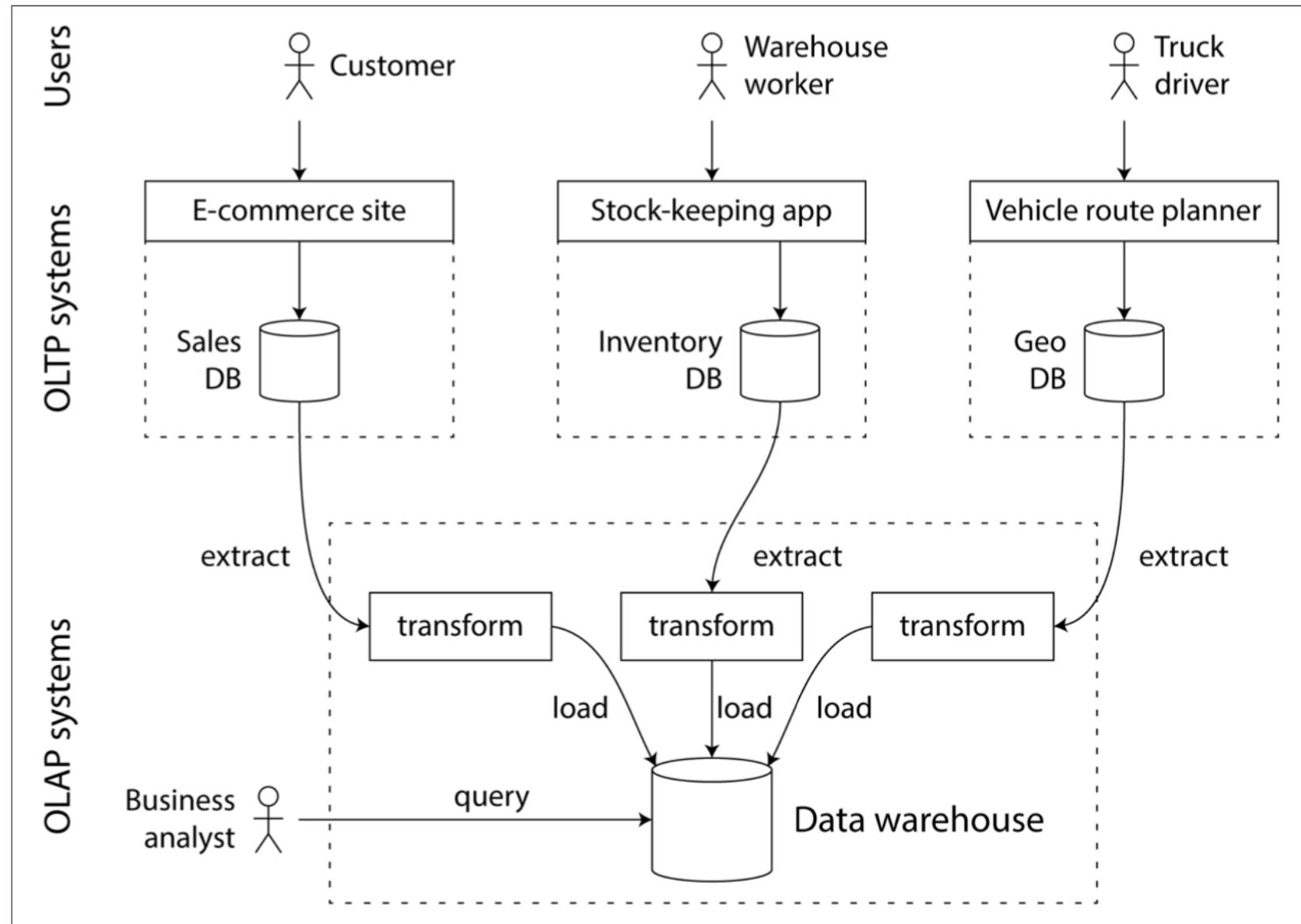
# What kind of data access do you need?

- **relational**: pandas, SQL: Postgres, sqlite, Hbase, VoltDB
- **document oriented**: MongoDB, CouchDB
- **key-value**: Riak, Redis, Memcached
- **graph oriented**: Neo4J

# What kind of data access do you need?

- in memory
- on disk
- on cluster

# OLTP vs OLAP



# Normalization to ....

