

DuckDB:

An embedded database for data science

Pedro Holanda

pedro@duckdblabs.com

- ▶ CWI in Amsterdam;  
- ▶ Database Architectures Group
- ▶ Why a new system?
 - ▶ Most Projects at CWI are sponsored by private companies;
 - ▶ Most of our projects are related to **data** science (Particularly, ML pipelines);
 - ▶ Main goal is to optimize/facilitate the management of data in data science projects.



Outline

- ▶ Why should I use a Database System?
- ▶ Combining Database Systems with Data Science.
- ▶ DuckDB: An embedded database system for data science.
- ▶ **Hands-on ~ 45 min.**
[Using DuckDB]

Why should I use a database system?

- ▶ Database that models a digital music store to keep track of artists and albums.
- ▶ Things we need to store:
 - ▶ Information about artists.
 - ▶ What albums those artists released.

- ▶ Store database as comma-separated value (CSV) files that we manage in our own code
 - ▶ Use separate file per entity
 - ▶ The application has to parse files each time they want to read/update records

- ▶ Database that models a digital music store

Artist (name, year, country)

"Backstreet Boys", 1994, "USA"

"Ice Cube", 1992, "USA"

"Notorious BIG", 1989, USA

Album (name, artist, year)

"Millenium", "Backstreet Boys", 1999

"DNA", "Backstreet Boys", 2019

"AmeriKKKa's Most Wanted", "Ice Cube", 1990

- ▶ Get the year that Ice Cube went solo

Artist (name, year, country)

"Backstreet Boys",1994,"USA"

"Ice Cube", 1992,"USA"

"Notorious BIG",1989,USA



```
1  for line in file
2      record = parse(line)
3      if 'Ice Cube' == record[0]
4          print int(record[1])
5
```

- ▶ How do we ensure that the artist is the same for each album entry?
- ▶ What if someone overwrites the album year with an invalid string?
- ▶ How do we store that there are multiple artists on an album?

- ▶ How do we find a particular record?
- ▶ What if we now want to create a new application that uses the same database?
- ▶ What if two threads try to write to the same file at the same time?

- ▶ Software that allows application to store and analyse information in a database.
- ▶ A general-purpose DBMS is designed to allow the definition, creation, querying, update and administration of databases.



- ▶ Many data scientists do not use database systems
 - ▶ Despite requiring many of the things they offer!
 - ▶ Data management, data wrangling...
- ▶ Instead: Engineer their own solutions
 - ▶ **Flat files** to store data (CSV, Binary, HDF5, etc).
 - ▶ **dplyr/pandas** as query execution engines.

- ▶ Manually managing files is cumbersome.
- ▶ Loading and parsing e.g. CSV files is inefficient.
- ▶ File writers typically do not offer resiliency.
 - ▶ Files can be corrupted;
 - ▶ Difficult to change/update.
- ▶ It does not scale!
- ▶ Why people use it?
 - ▶ A pip install and you ready to go.

```
$ pip install pandas, numpy
```

- ▶ The problem is that they are *very poor query engines!*
- ▶ Materialize huge intermediates
- ▶ **No query optimizer**
 - ▶ Not even for basics like filter pushdown
- ▶ No support for out of memory computation
- ▶ No support for parallelization
- ▶ Unoptimized implementations for joins/aggregations

- ▶ Data scientists **need** the functionality database systems offer
- ▶ But they opt not to use database systems
- ▶ Often this leads to problems down the road
 - ▶ When the data gets bigger...
 - ▶ When a power outage corrupts their data...

Combining Database Systems with Data Science

- ▶ How can we combine analytical tools (R/Python) with database systems?

- ▶ Database Client Connections
- ▶ User Defined Functions
- ▶ Embedded Database Systems

▶ 1: DB Connection

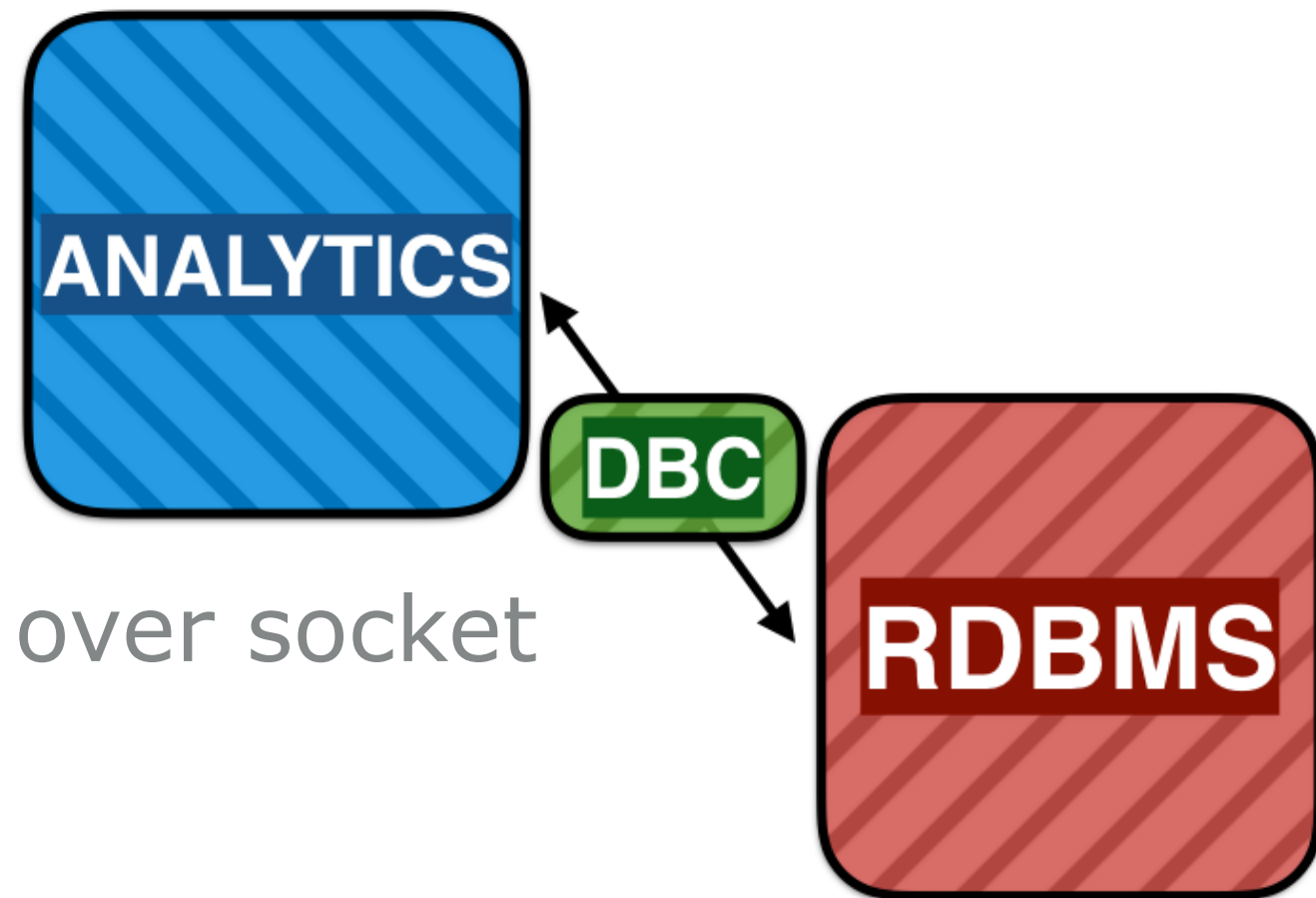
▶ DBMS is separate process

▶ Queries & Data transferred over socket

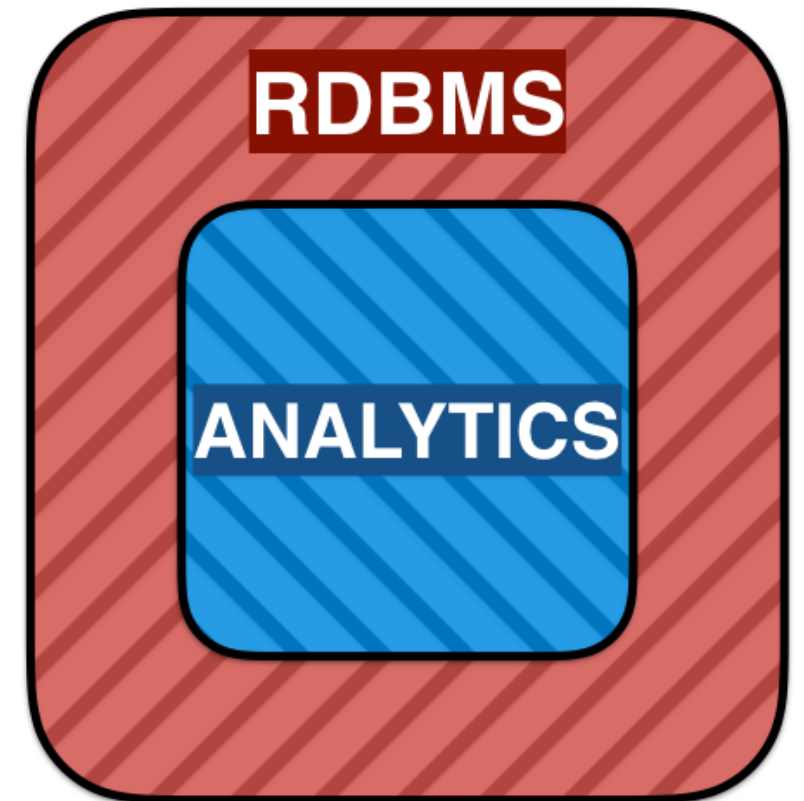
▶ Problems:

▶ Data transfer is very slow (both directions)

▶ Requires setup & management of DBMS server



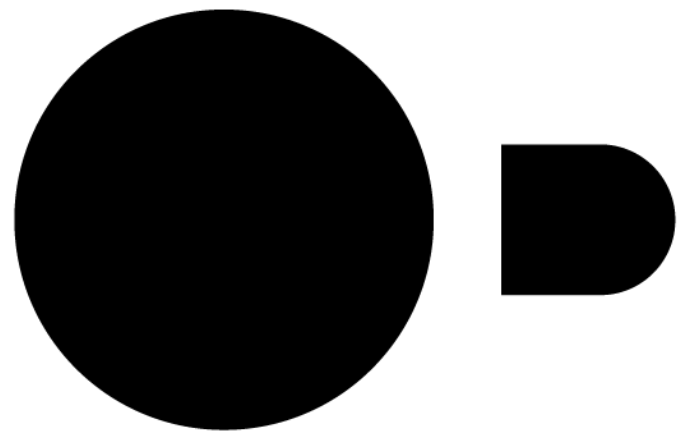
- ▶ Most cases a data scientist exports data from the RDBMS to the analytical tool
- ▶ **User-Defined Functions (UDFs)**
- ▶ Analytics is run inside DBMS server
- ▶ No separate analytics program!
- ▶ Problems:
 - ▶ Difficult to implement and debug
 - ▶ DBMS-specific, requires knowledge of DBMS internals
 - ▶ Also requires setup & management of DBMS server



- ▶ It runs inside the analytics applications;
- ▶ Has same low-transfer cost advantages as UDFs;
- ▶ Are easy to install/use;
- ▶ Binds to almost every language.
- ▶ What the most famous embeddable DBMS?
 - ▶ Runs on every phone, browser and OS
 - ▶ It even runs inside airplanes!
- ▶ **Great for transactions, not so good for analytics.**



CWI



DuckDB

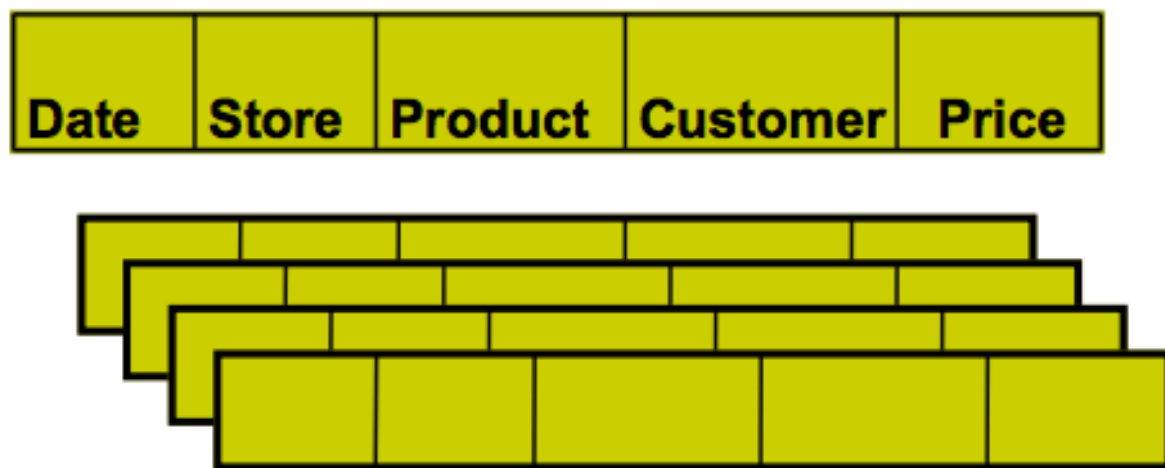
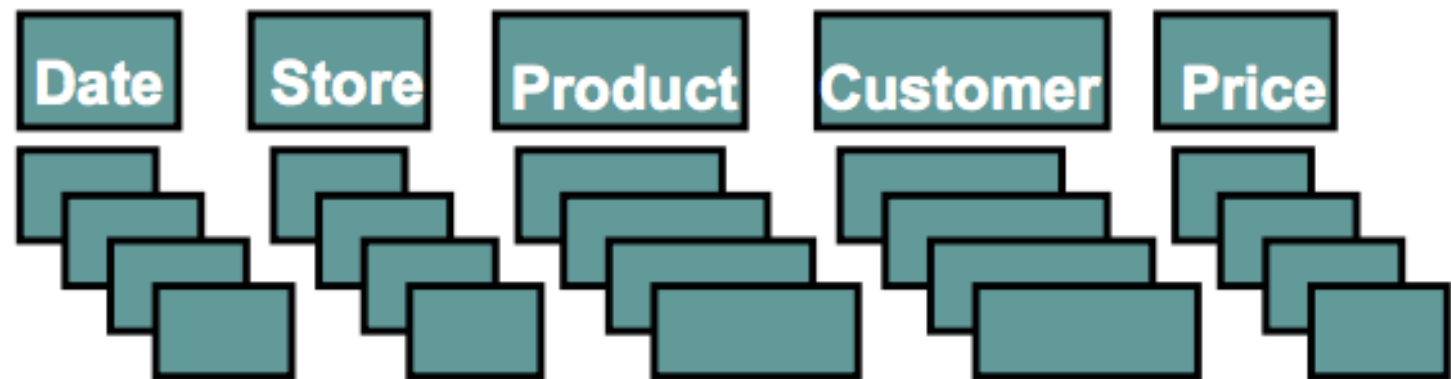
- ▶ DuckDB: The SQLite for Analytics
- ▶ **Simple installation**

```
$ pip install duckdb
```
- ▶ **Embedded: no server management**
- ▶ **Fast analytical processing**
- ▶ **Fast transfer between R/Python and RDBMS**
- ▶ **Duckdb is currently in pre-release**
 - ▶ Check duckdb.org for more details.



- ▶ Data Science is equal to Analytical Processing!
- ▶ Storage Model
 - ▶ Row-Store vs Column-Store
- ▶ Compression
- ▶ Query Execution
 - ▶ Row-wise vs vector-wise
- ▶ Relational API

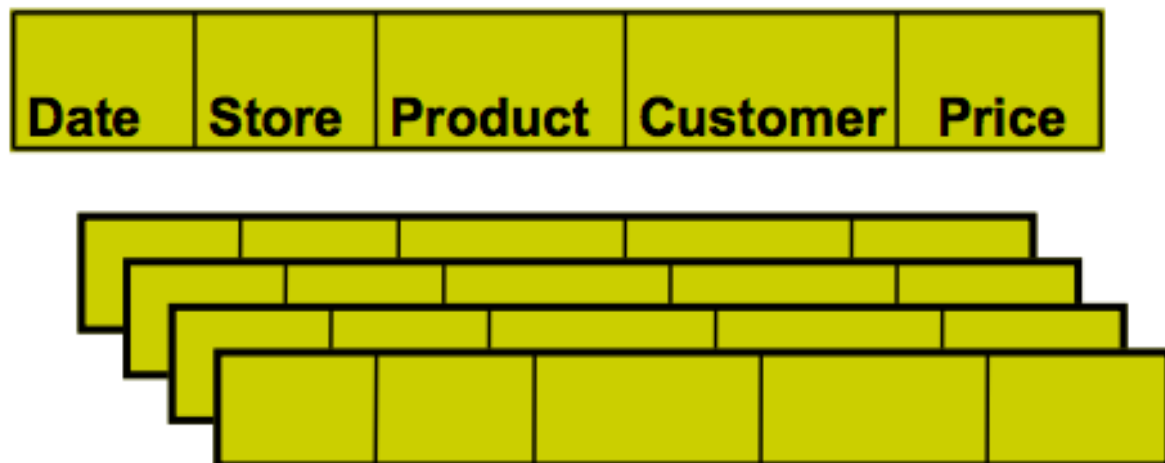
- ▶ SQLite use a row-storage model
- ▶ DuckDB uses a columnar storage model

row-store**column-store**

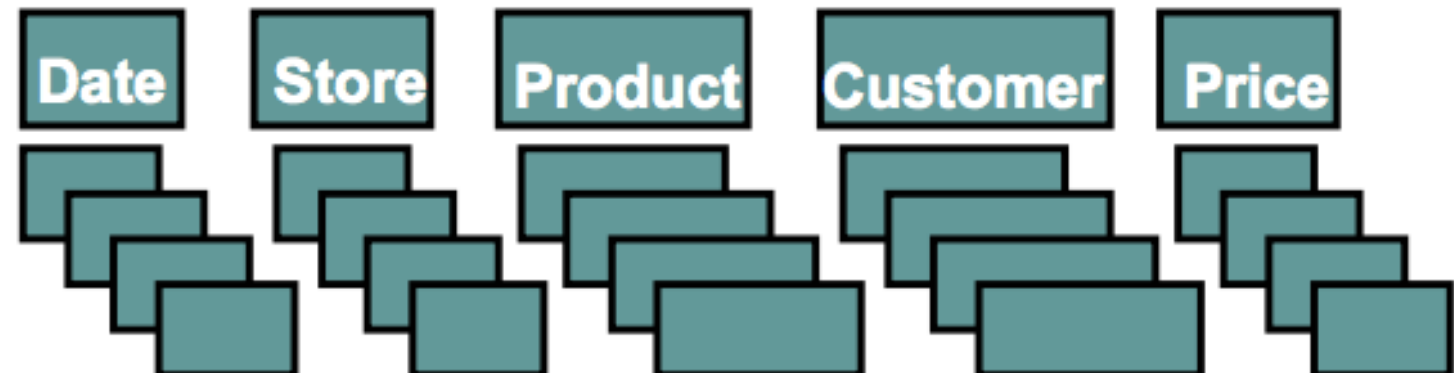
▶ Row-Storage:

- ▶ Individual rows can be fetched cheaply
- ▶ However, **all columns must always be fetched!**
- ▶ What if we only use a few columns?
- ▶ **e.g.:** What if we are only interested in the price of a product, not the stores in which it is sold?

row-store



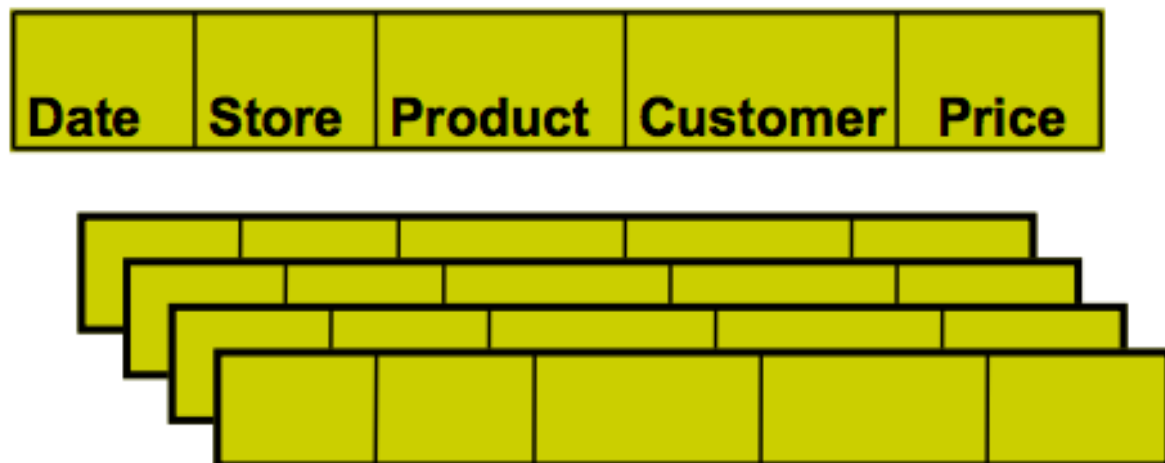
column-store



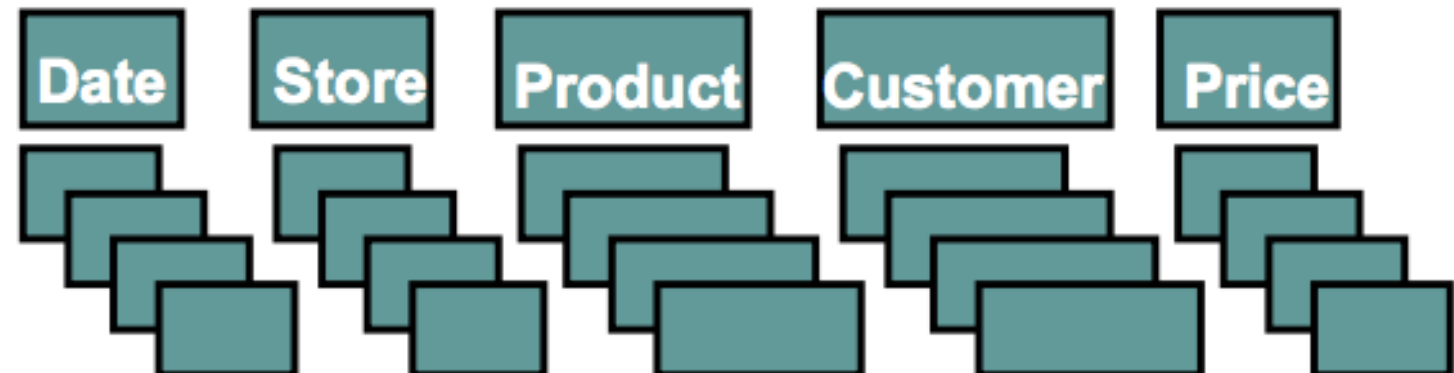
▶ Column-Storage:

- ▶ We can fetch individual columns
- ▶ Immense savings on disk IO/memory bandwidth when only using few columns

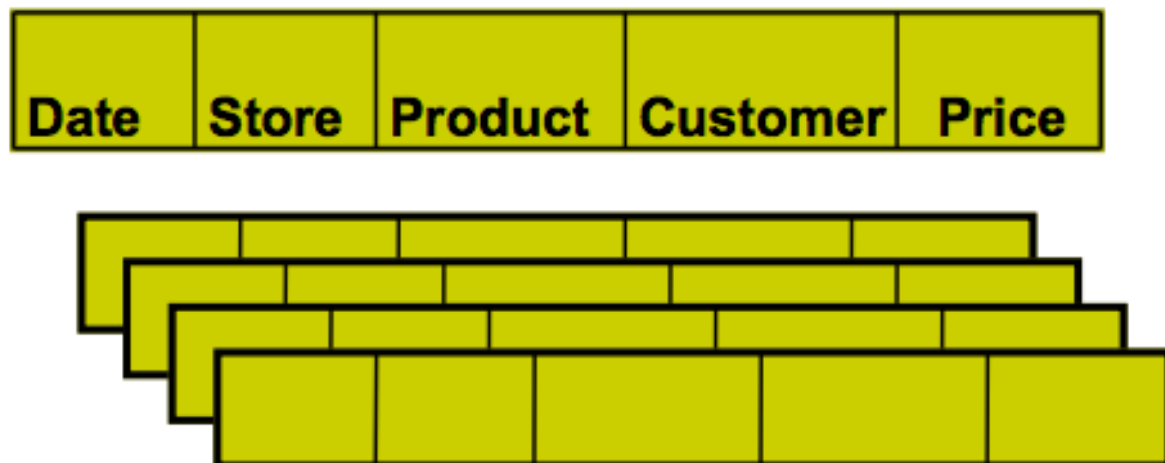
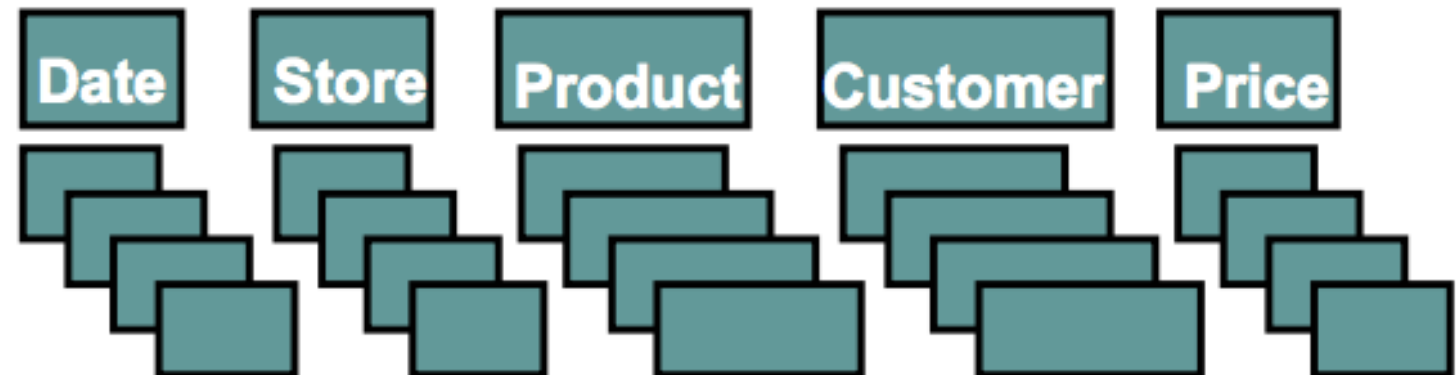
row-store



column-store

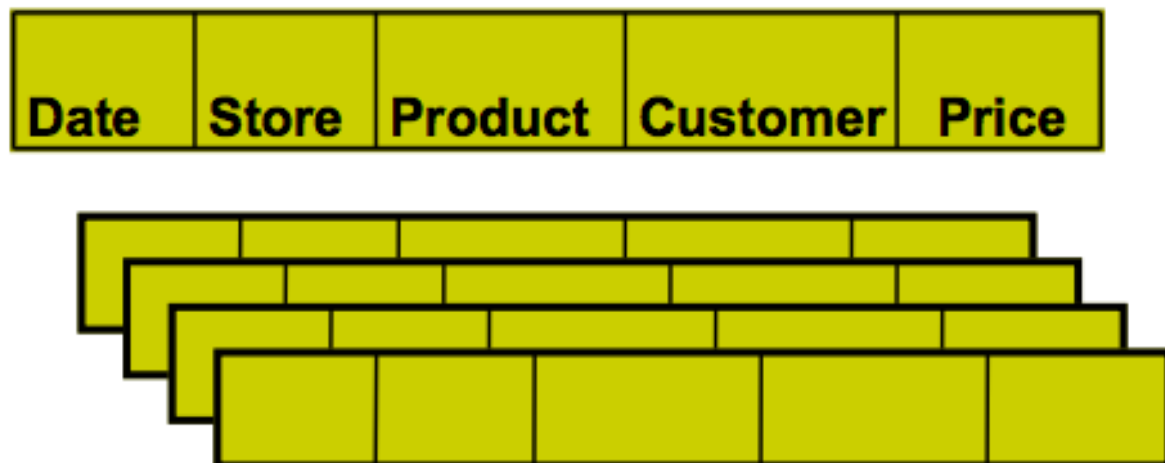


- ▶ **Example:** Suppose we have a 1TB table with 100 columns. We have a query that requires 5 columns of the table.
- ▶ **Row-store:** Read entire 1TB of data from disk at 100MB/s \cong 3 hours
- ▶ **Column-store:** Read 5 columns (50GB) from disk \cong 8 minutes

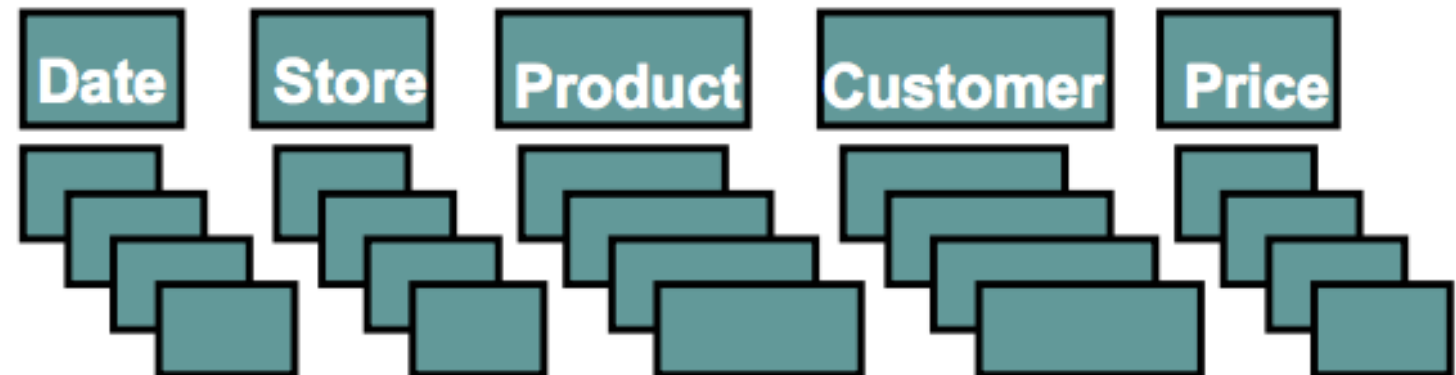
row-store**column-store**

- ▶ **Compressibility** is another advantage of column-storage
- ▶ Individual columns often have similar values, e.g. dates are usually increasing
- ▶ Save ~**2-10X** on storage (depending on compression algorithms used and data)

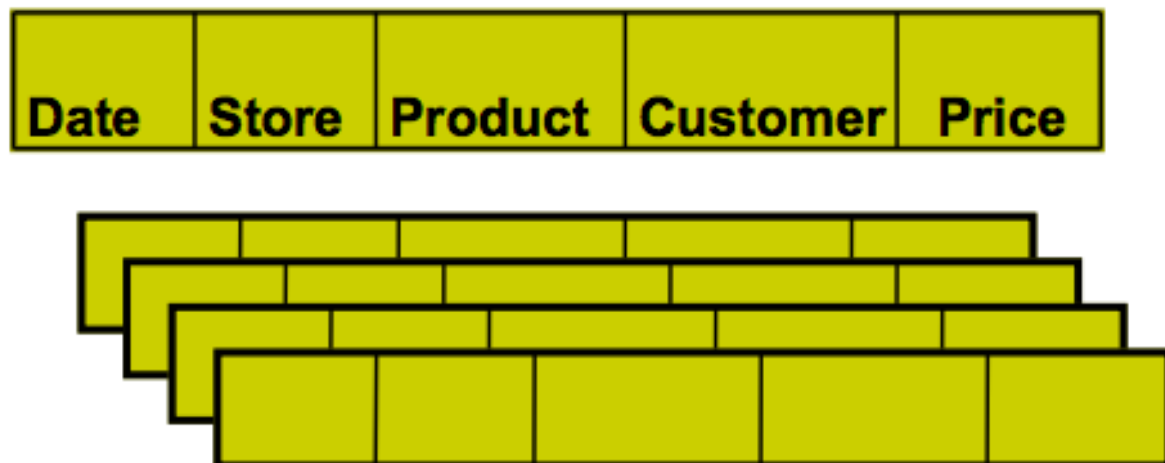
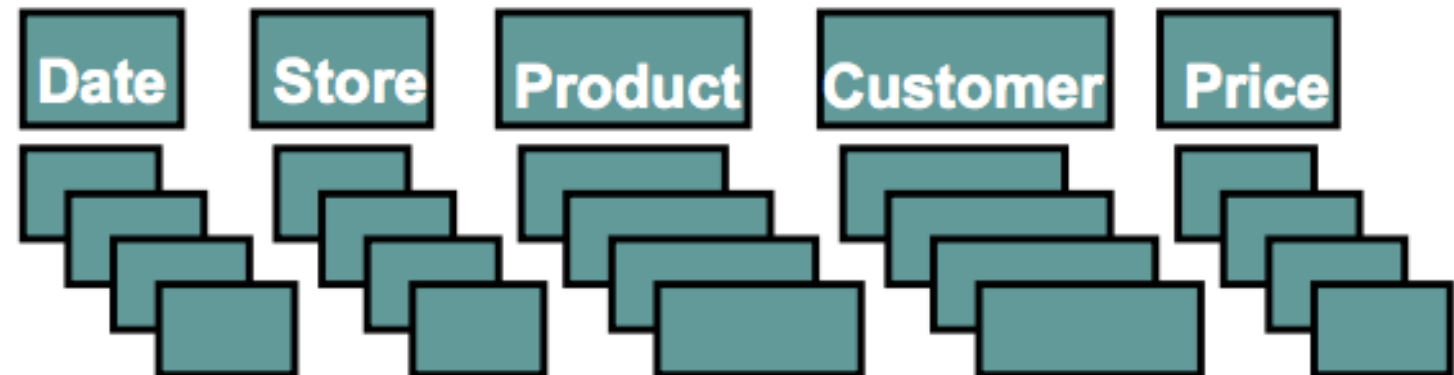
row-store



column-store

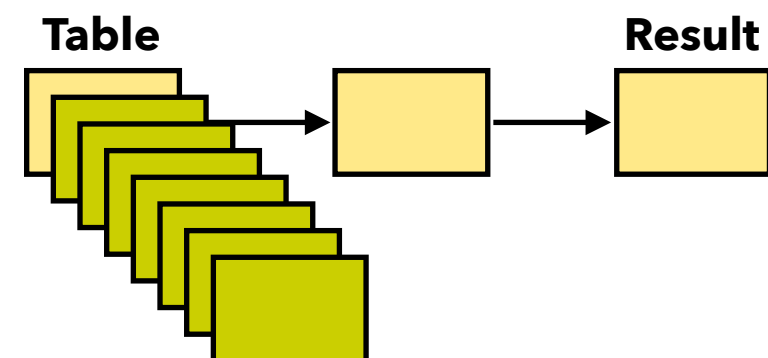


- ▶ **Example:** Suppose we have a 1TB table with 100 columns. We have a query that requires 5 columns of the table.
- ▶ **No compression:** Read 5 columns (50GB) from disk \approx 8 minutes
- ▶ **Compression:** Read 5 compressed columns (5GB) from disk \approx 50 seconds

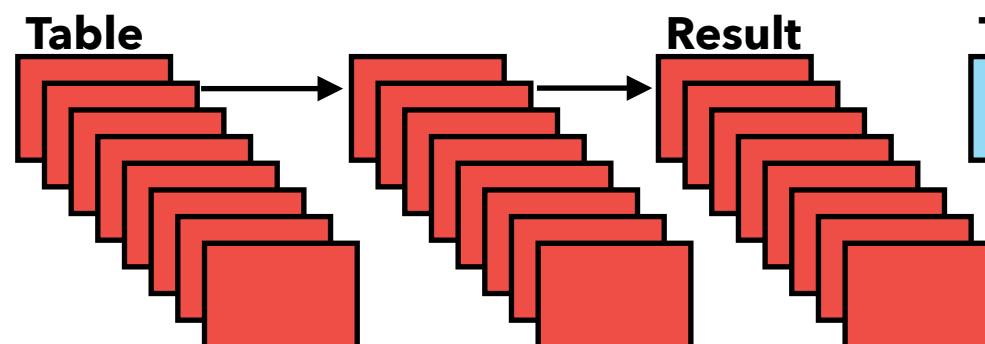
row-store**column-store**

- ▶ **SQLite** use tuple-at-a-time processing
 - ▶ Process **one row** at a time
- ▶ **NumPy/R** use column-at-a-time processing
 - ▶ Process entire columns at once
- ▶ **DuckDB** uses **vectorized** processing
 - ▶ Process **batches** of columns at a time

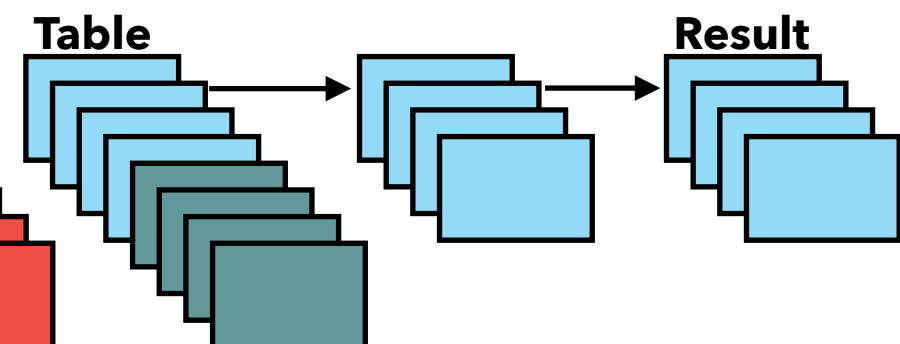
Tuple-at-a-Time



Column-at-a-Time

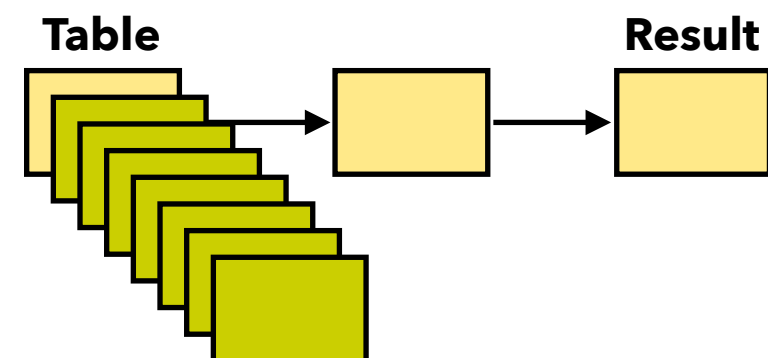


Vectorized Processing

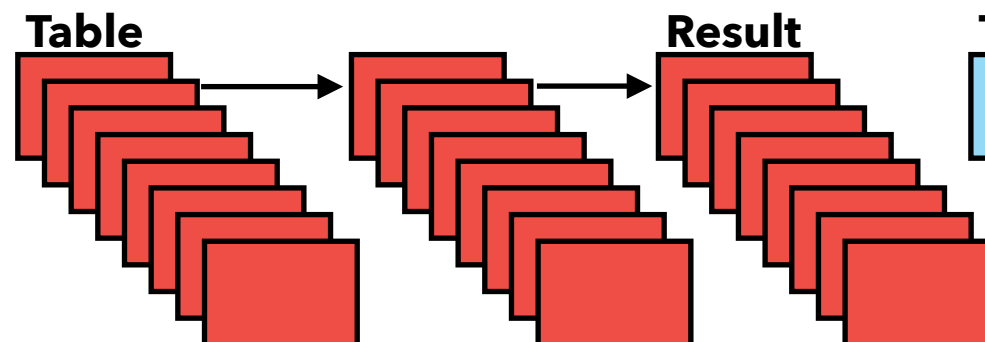


- ▶ **Tuple-at-a-Time (SQLite)**
 - ▶ Optimize for low memory footprint
 - ▶ Only need to keep **single row** in memory
- ▶ Comes from a time when **memory was expensive**
- ▶ **High CPU overhead per tuple!**

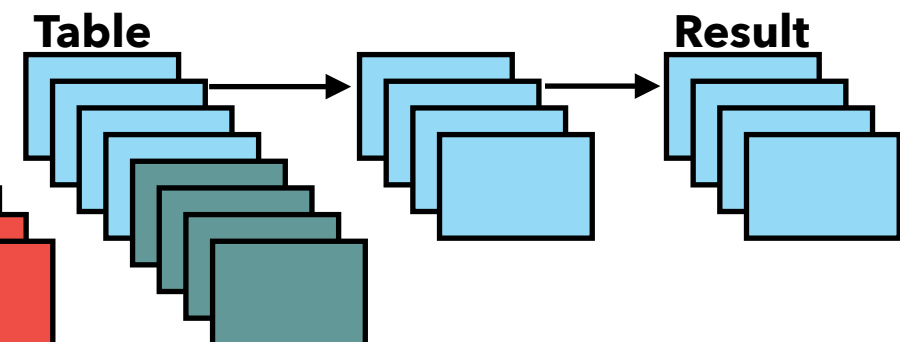
Tuple-at-a-Time



Column-at-a-Time

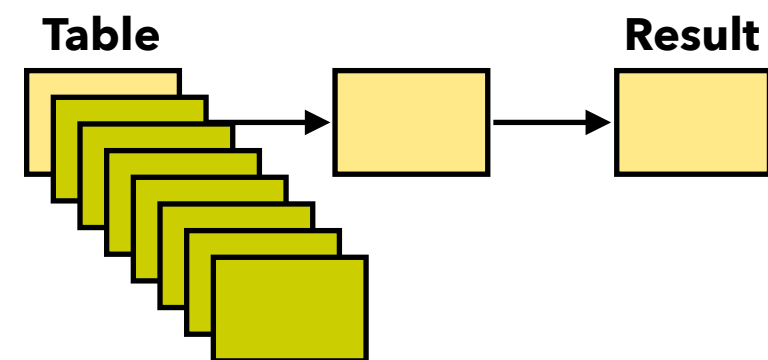


Vectorized Processing

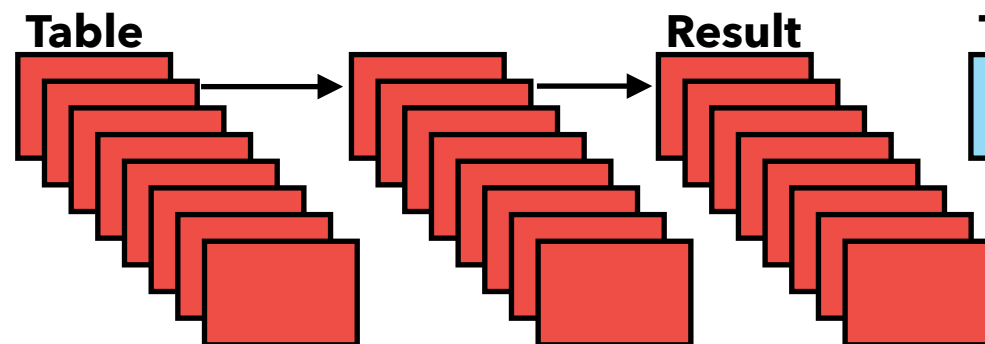


- ▶ **Column-at-a-Time (NumPy/R)**
 - ▶ Better CPU utilization, allows for SIMD
 - ▶ Materialize **large intermediates** in memory!
- ▶ Intermediates can be gigabytes each...
- ▶ **Problematic** when data sizes are large

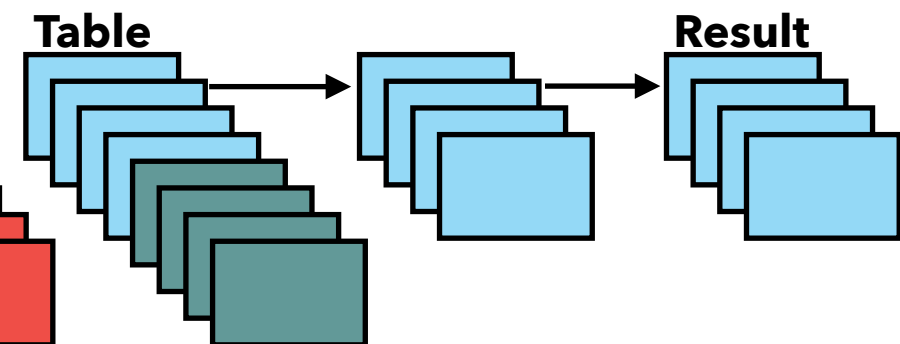
Tuple-at-a-Time



Column-at-a-Time

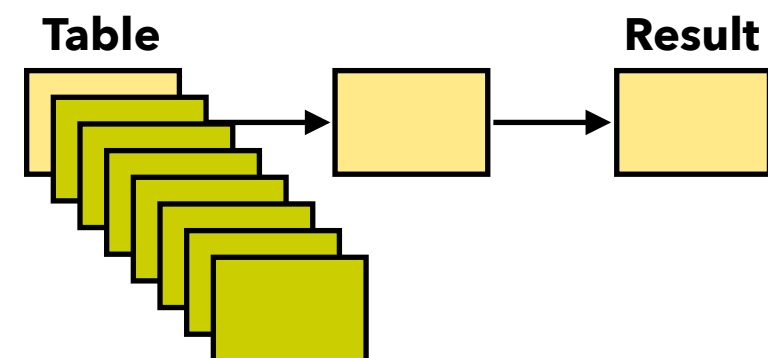


Vectorized Processing

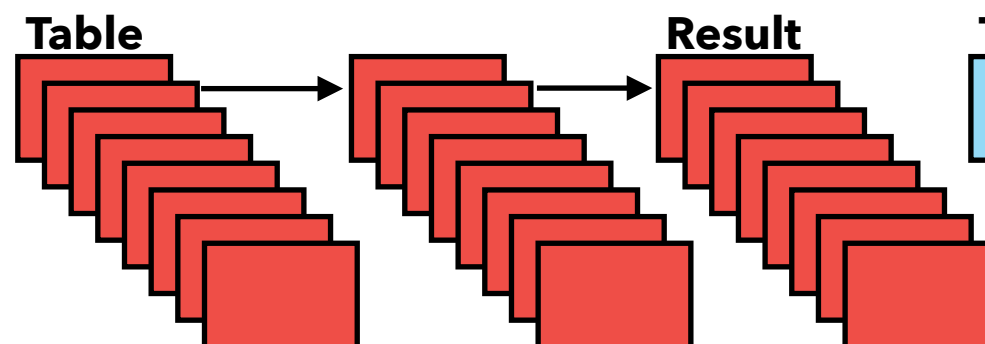


- ▶ **Vectorized Processing (DuckDB)**
 - ▶ Optimized for CPU Cache locality
 - ▶ SIMD instructions, Pipelining
 - ▶ **Small intermediates (ideally fit in L1 cache)**

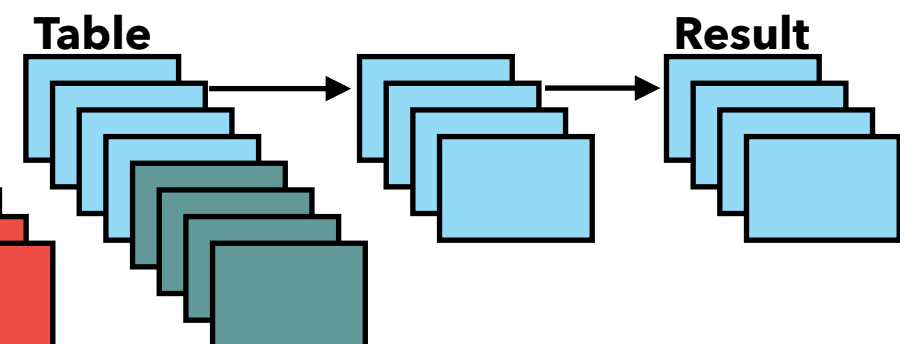
Tuple-at-a-Time



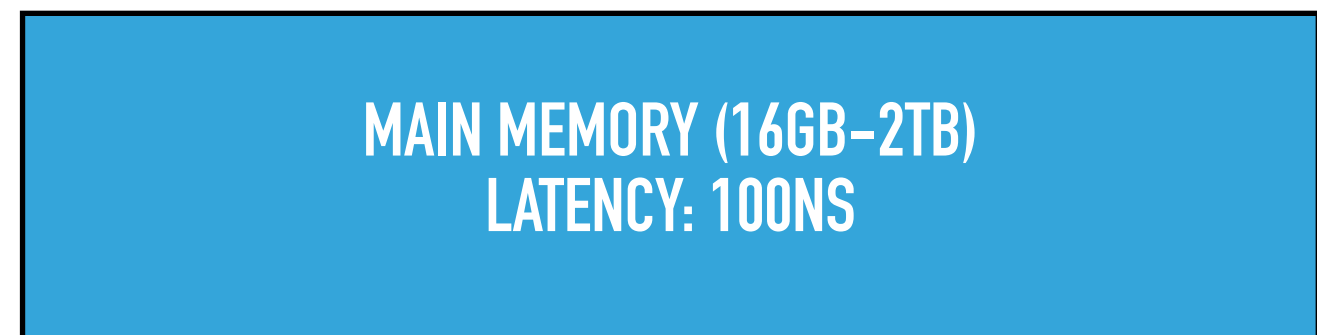
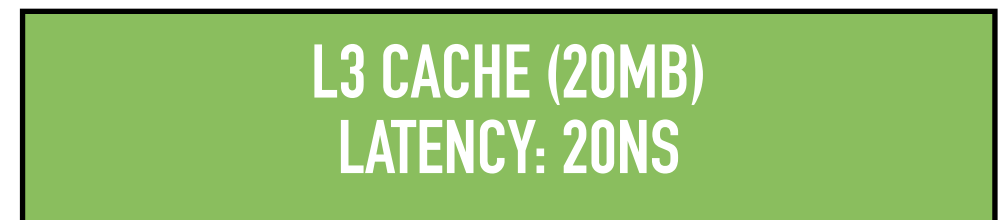
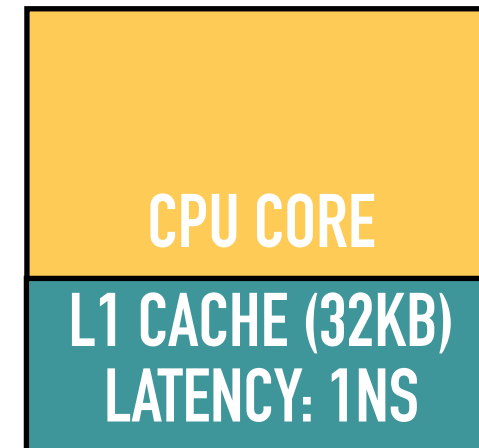
Column-at-a-Time



Vectorized Processing



- ▶ **Vectorized Processing**
- ▶ Intermediates fit in L3 cache
- ▶ **Column-at-a-Time**
- ▶ Intermediates go to memory



CPU CORE

**L1 CACHE (32KB)
LATENCY: 1NS**

**L2 CACHE (256KB)
LATENCY: 5NS**

**L3 CACHE (20MB)
LATENCY: 20NS**

**MAIN MEMORY (16GB-2TB)
LATENCY: 100NS**

- ▶ One of the reasons data scientists love libraries like pandas is their easy API.
- ▶ Database Systems forces developers to use SQL
- ▶ In DuckDB we have been developing a new API, called relational API that resembles the API of libraries like Pandas.

DB API

```
1 import duckdb
2 con = duckdb.connect('database.db')
3 cursor = db.cursor()
4 cursor.execute('Select j+1 from integers where i = 2')
```

RELATIONAL API

```
1 import duckdb
2 con = duckdb.connect('database.db')
3
4 # table operator returns a table scan
5 tbl = duckdb.table('integers')
6 # we can inspect intermediates
7 tbl.show()
8 # we can chain multiple operators
9 tbl.filter('i=2').project('j+1').show()
```

Hands-on

- ▶ Goal: See in practice the differences of:
 - ▶ Pandas;
 - ▶ DuckDB.
- ▶ Tasks (Benchmark):
 - ▶ Queries (e.g., aggregations, filters and joins);
 - ▶ Transactions.

- ▶ Right Now:
 - ▶ Clone Repo:
<https://github.com/pdet/duckdb-tutorial>
 - ▶ Upload file: Part 1/Exercise/Exercise.ipynb to
<https://colab.research.google.com/> as a Python 3 Notebook.