# Data Wrangling

## DuckDB

DuckDB is a lightweight OLAP type database system written in C++ and designed to be used for EDA style activities:



**From their website: advice on when to use and not to use DuckDB**

## Polars

Polars is a rust-based data frames library with Python bindings

[Here is a talk](#) that Juan Luis gave about the library

---