

NoDB: Efficient Query Processing on Raw Data Files

Ioannis Alagiannis, Renata Borovica, Miguel Branco, Stratos Idreos, Anastasia Ailamaki

In the Era of Data Deluge

Data collections become larger and larger

Many applications avoid using DBMS

e.g. social networks, scientific experiments

Problem: DBMS **startup cost**

- Why wait for loading and tuning?
- Why load the entire data set?
- Column Vs. row-store?
- Hire DB expert?

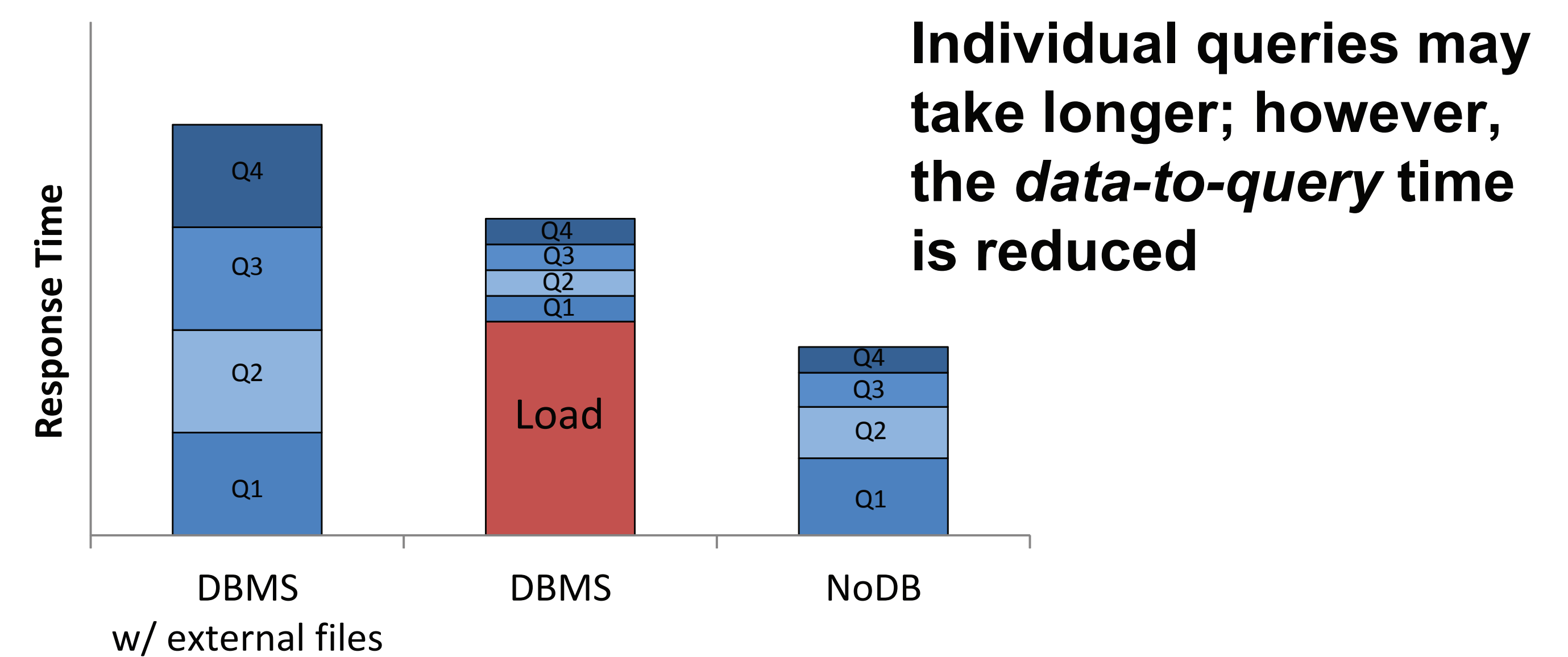


Challenges

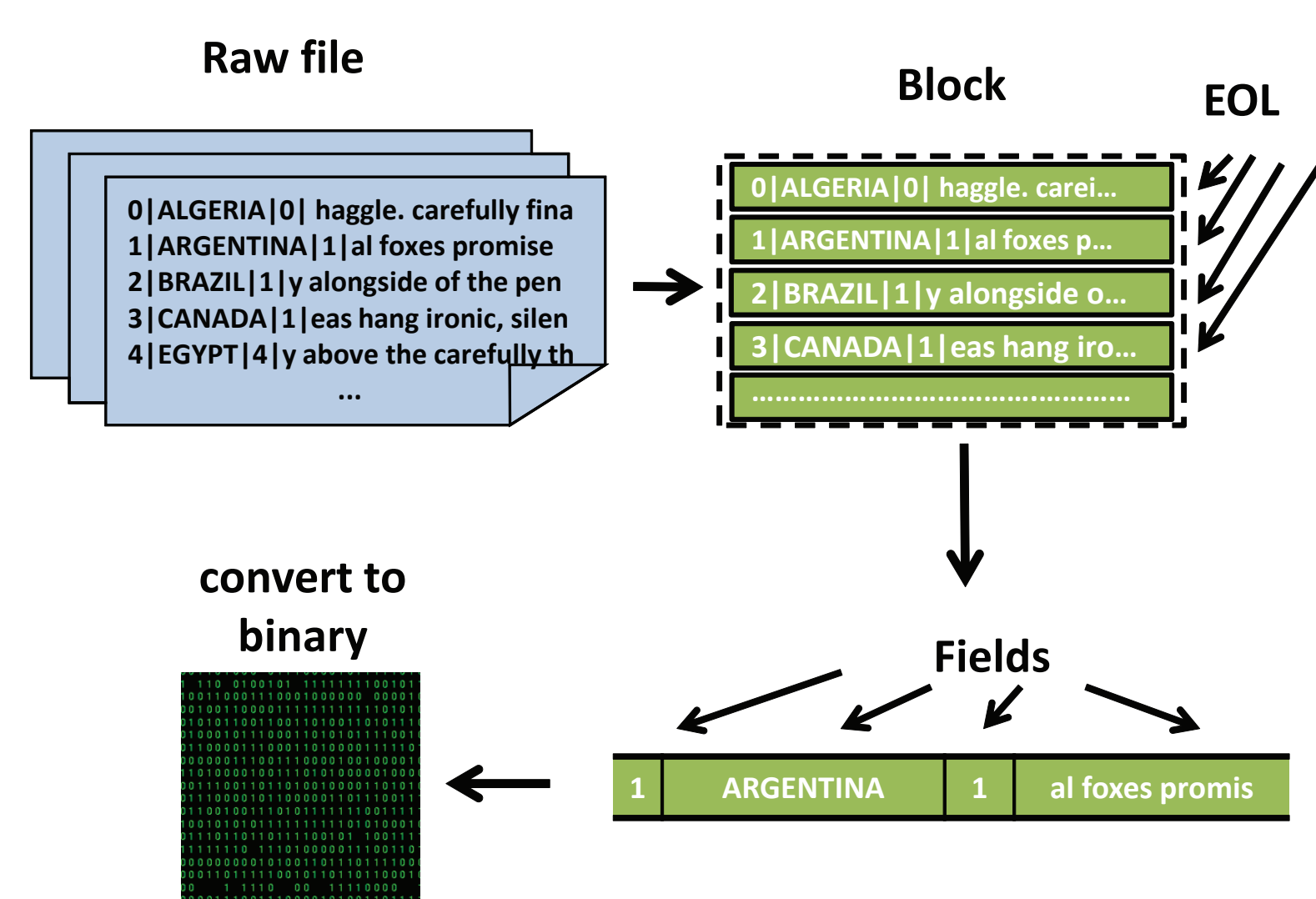
- Make DBMS more accessible to the user
- Eliminate major bottlenecks i.e. data loading

NoDB Philosophy

- No a priori data loading
- Raw data files as a first-class citizen
- Query processing *in situ*
- Maintain the whole feature set of a modern DBMS



Querying Raw Data Files



Positional Map

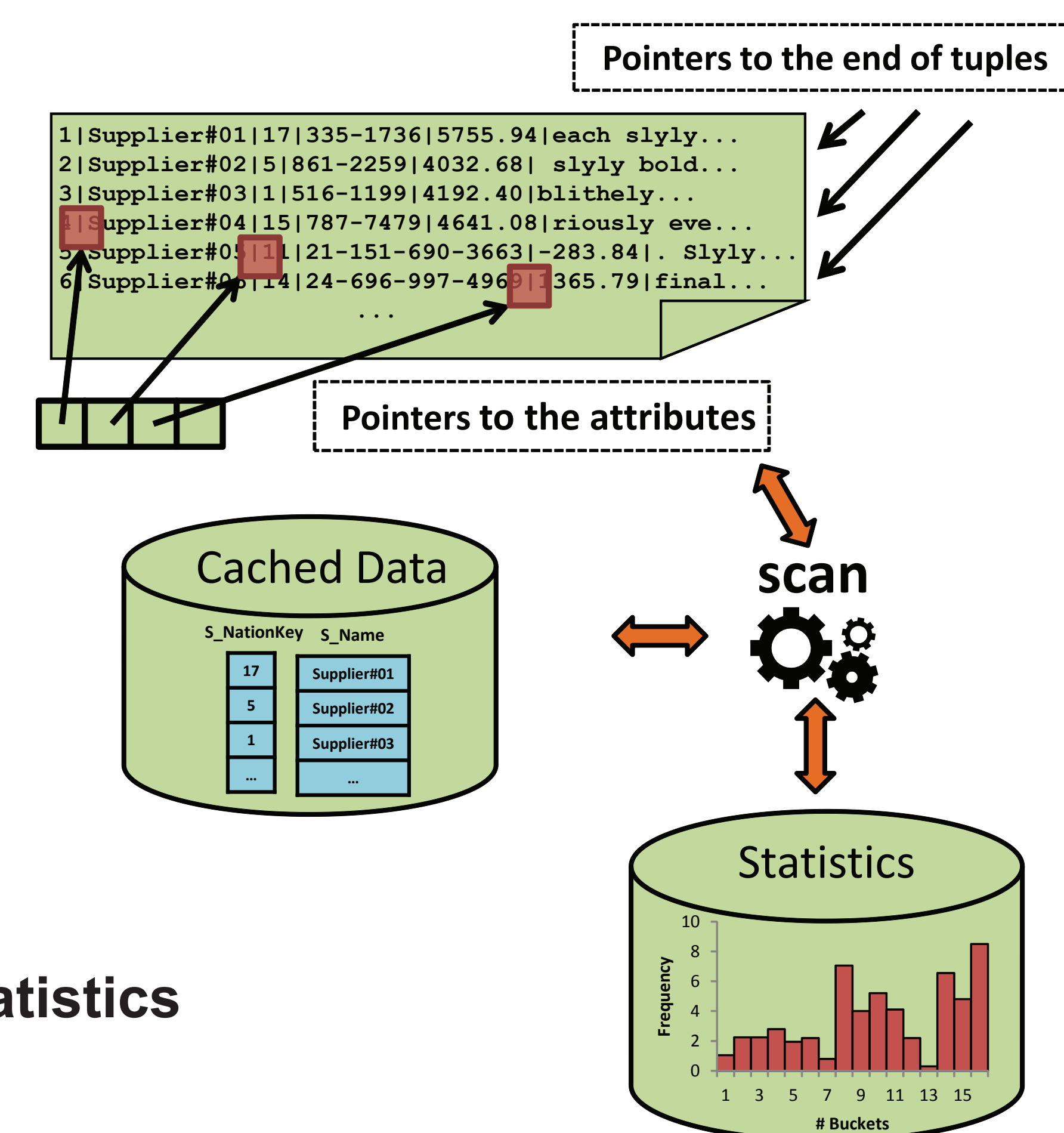
- Reduce parsing and tokenizing
- Maintain low level metadata
- Create on-the-fly
- Adapt to queries

Caching

- Complementary to positional map
- Avoid raw file accesses
- Populated during query execution

Statistics

- Extend the *scan* operator to create statistics
- Generate in an adaptive way



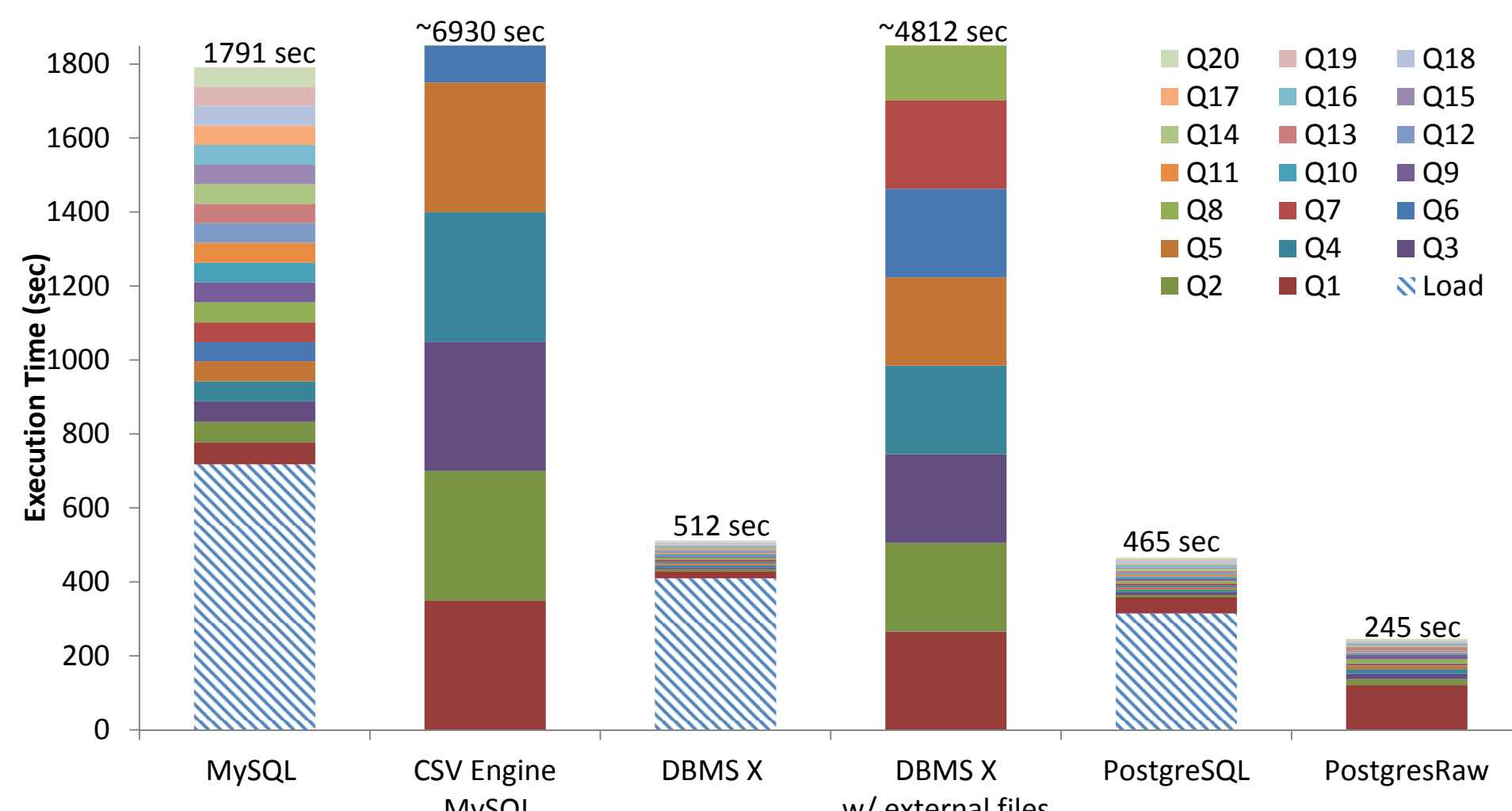
Accessing raw data files is expensive

- Parsing & tokenizing
- Data type conversion

PostgresRaw in Action

Evaluation

- Applying the NoDB philosophy to PostgreSQL without altering the internals of the query engine
- PostgresRaw vs. other DBMS



Micro Benchmark

- Data exploration scenario
- Tuples: 7.5 millions - Attributes: 150
- Epoch 1: access c1-c30
- Epoch 2: access c100-c130
- 20 aggregation queries in each epoch
- Vary selectivity

PostgresRaw is competitive with DBMS that have already loaded the data

Significant gains over external files

In situ query processing overhead is amortized across a sequence of queries

Adaptation to workload changes

- Exploring a new region of the file
- Adjust positional map and cache automatically

