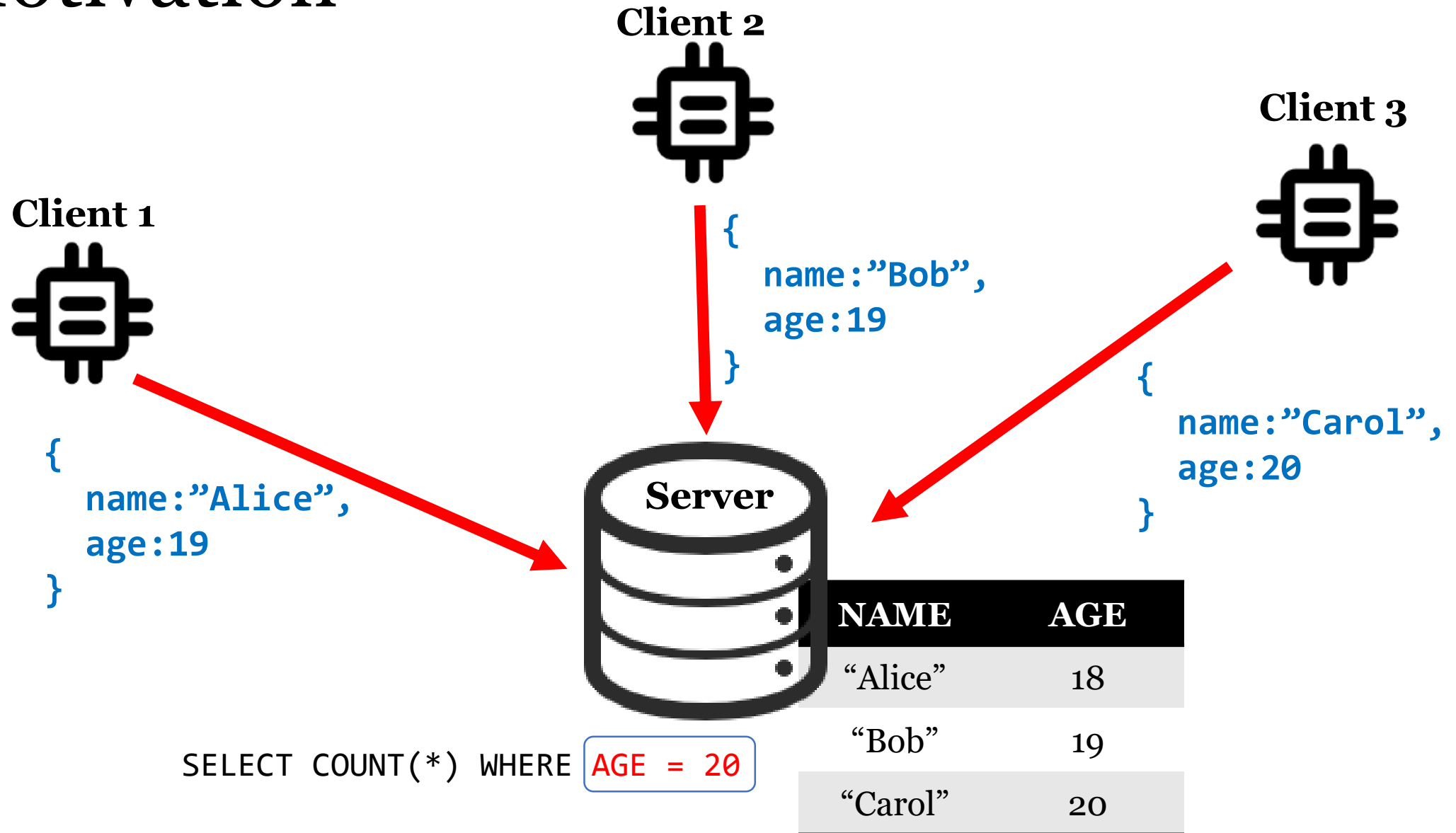


# CIAO: An Optimization Framework for Client-Assisted Data Loading

Cong Ding, Dixin Tang, Xi Liang, Aaron J. Elmore, Sanjay Krishnan

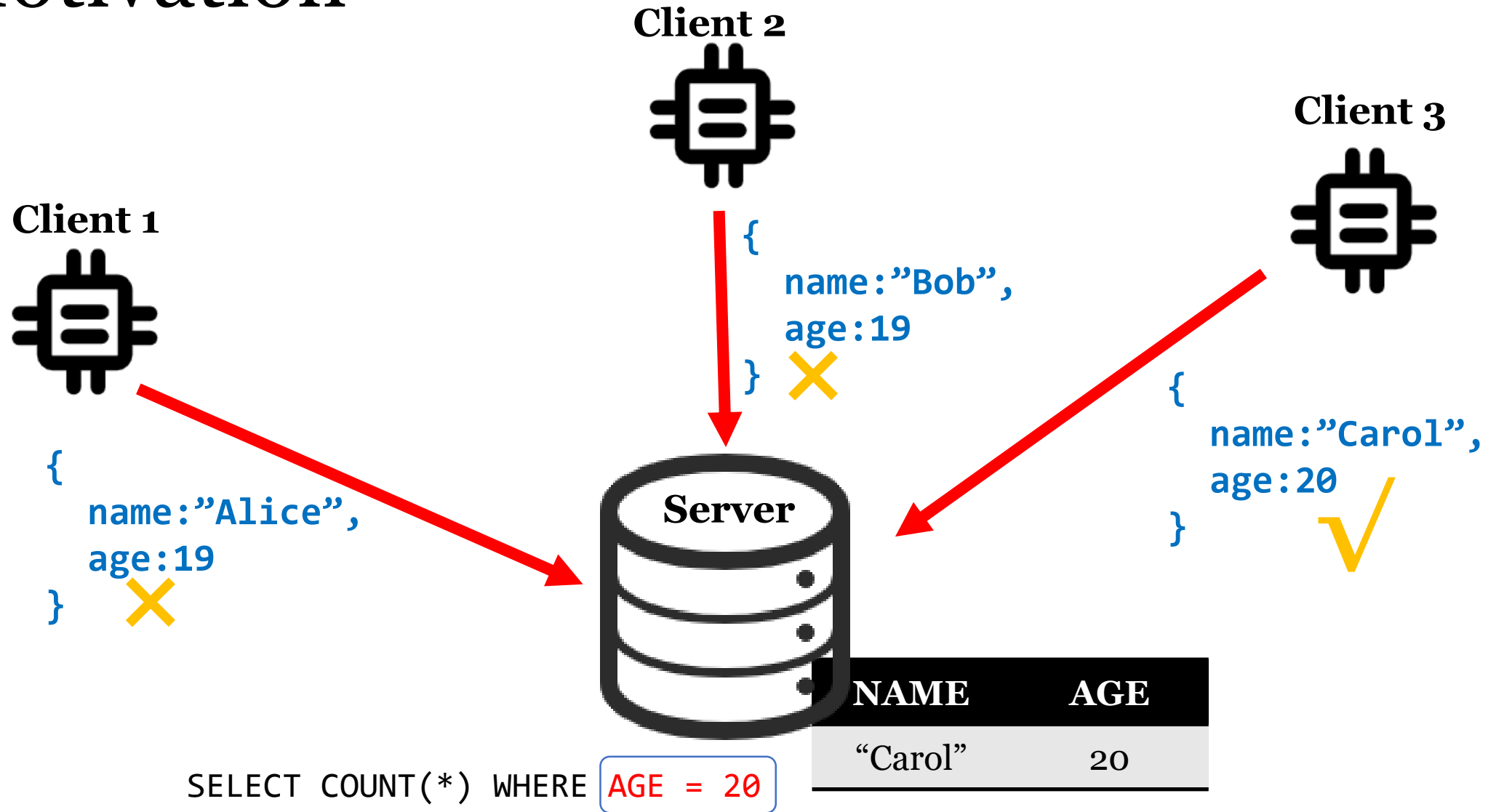


# Motivation

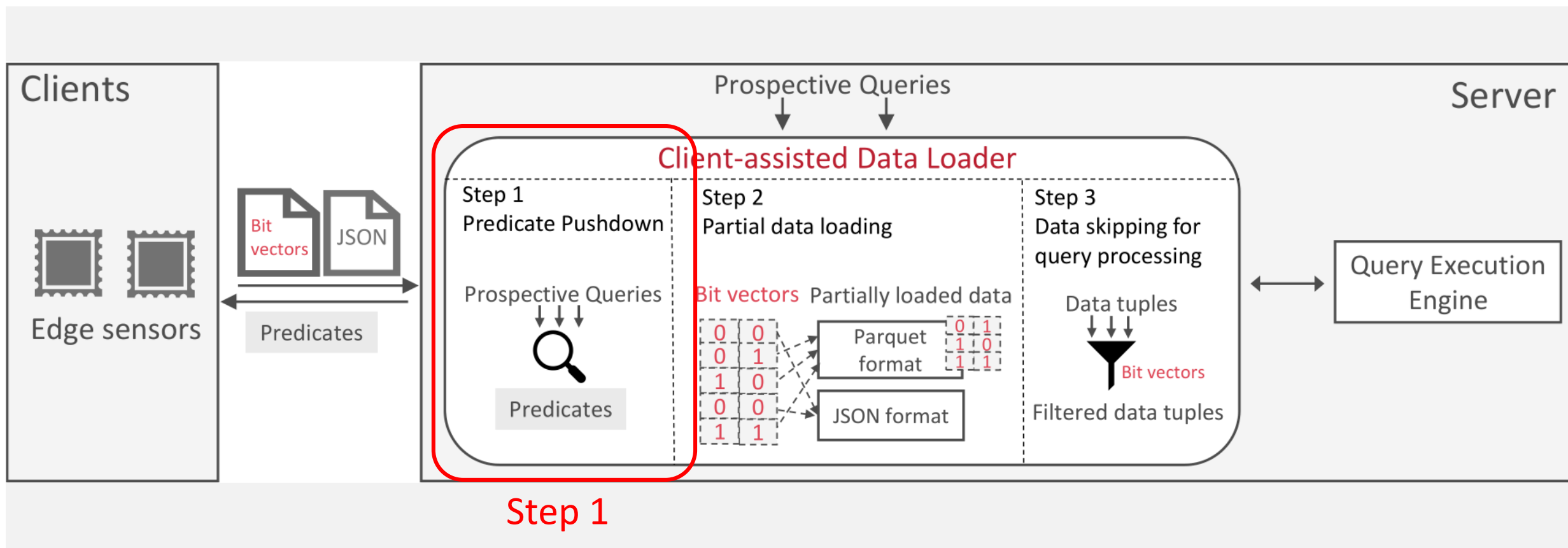


## Client-Assisted Data Loading

# Motivation



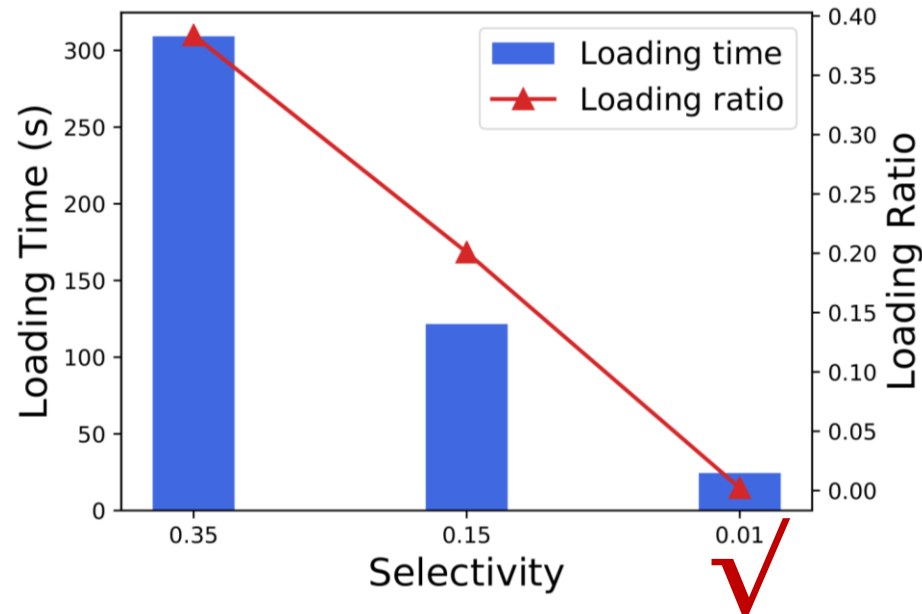
# Overview



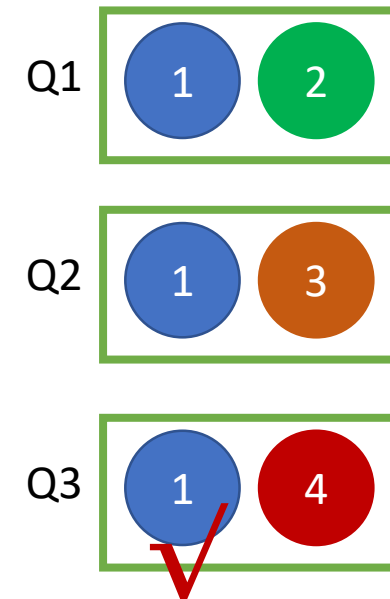
# Predicate Selection

- Evaluating predicate on the clients will incur computation cost
- Choose the most *beneficial* predicate set within *limited budget*

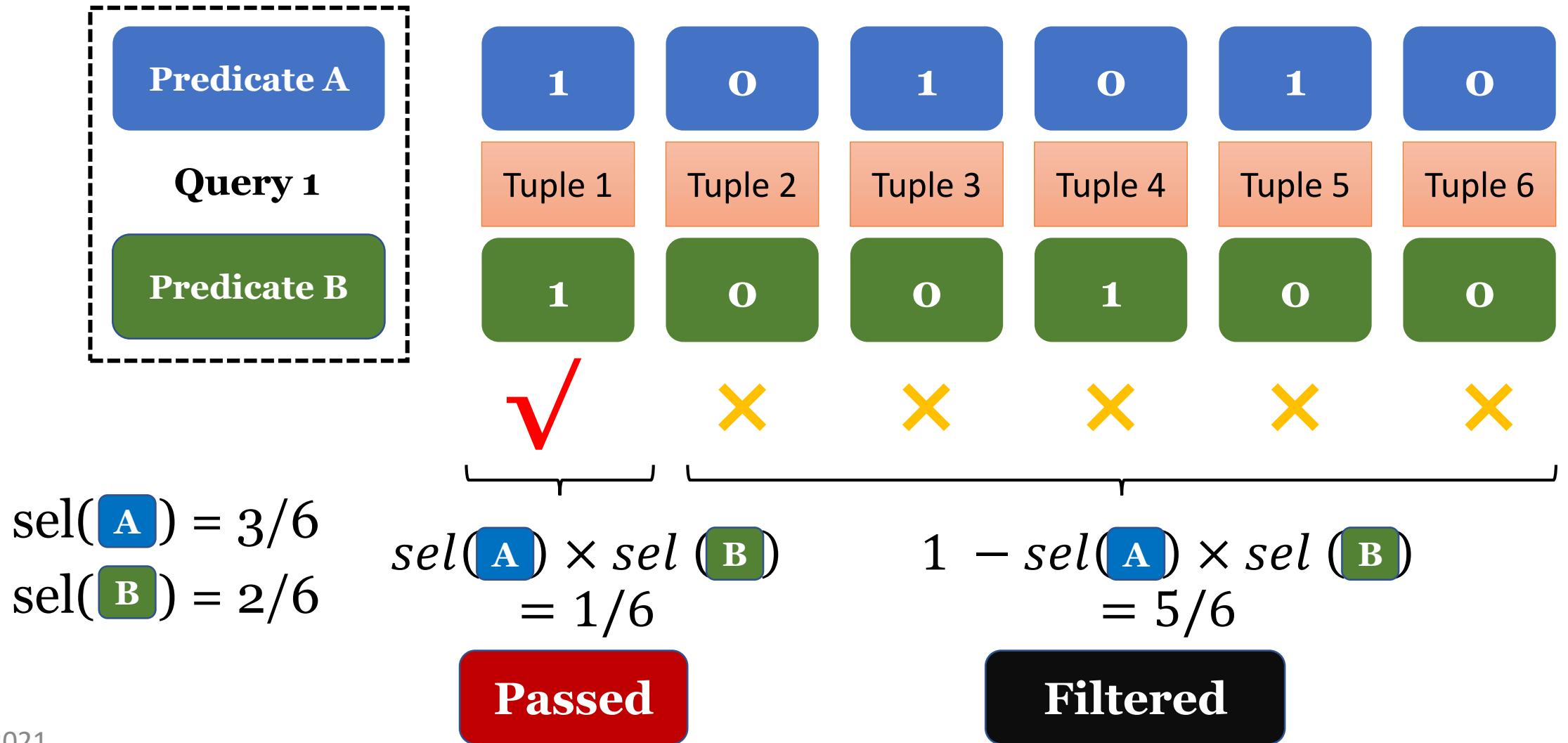
## Selectivity



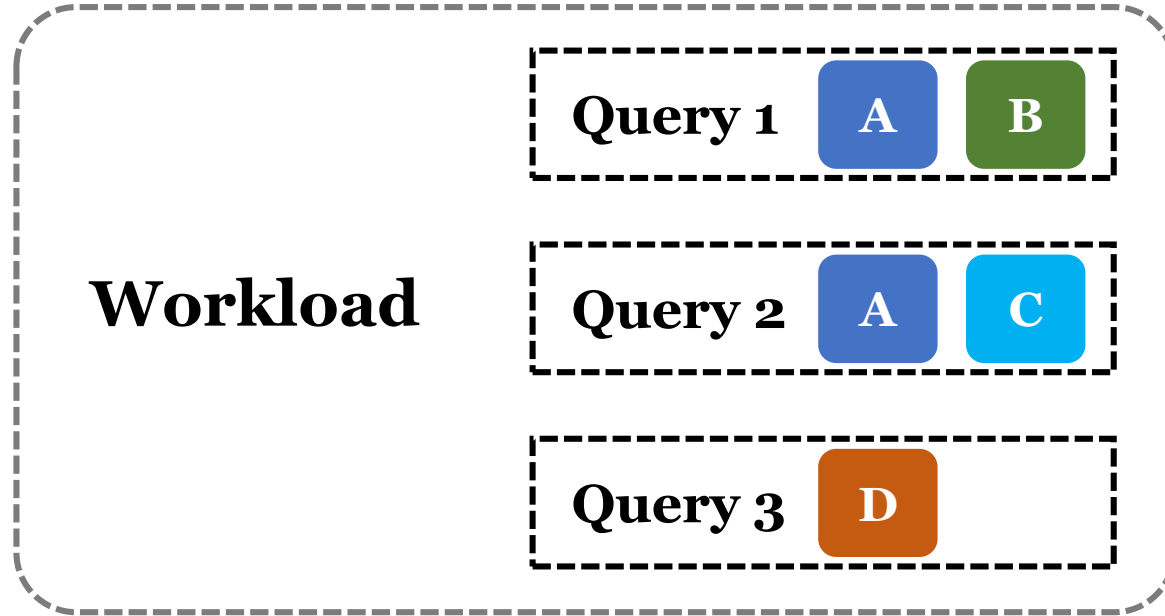
## Frequency



# Problem Formalization



# Problem Formalization



$$f = 1 - sel(\text{A}) \times sel(\text{B})$$

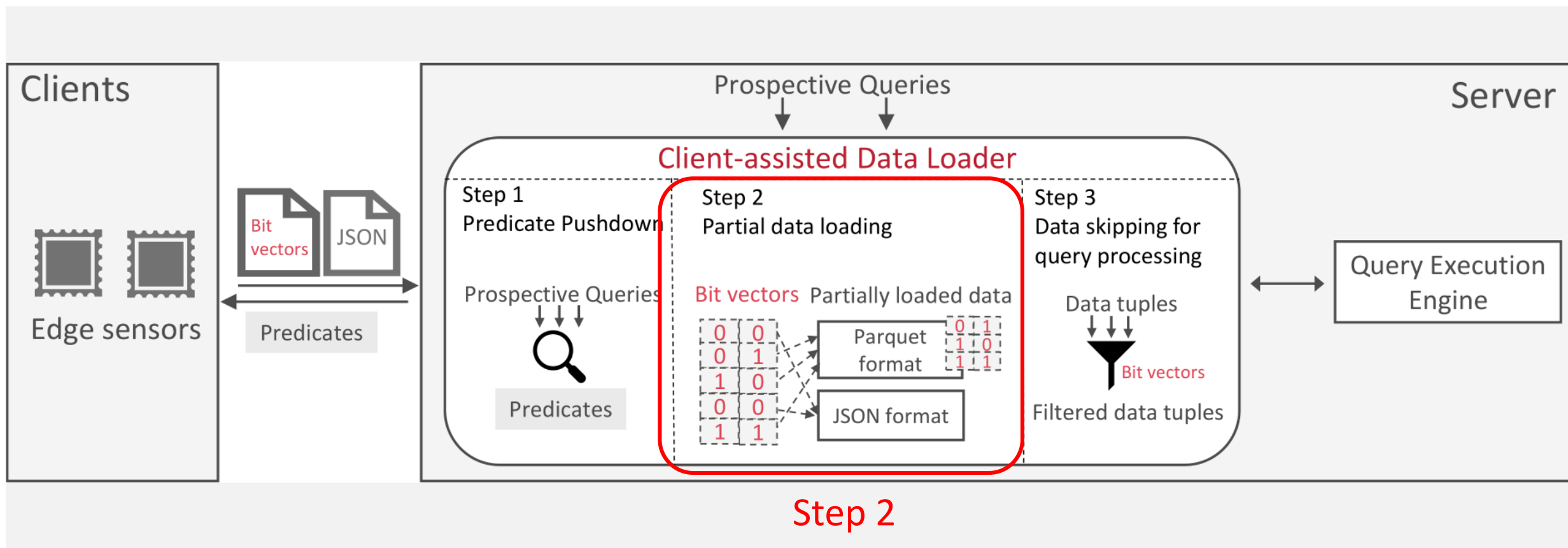
$$+ 1 - sel(\text{A}) \times sel(\text{C})$$

$$+ 1 - sel(\text{D})$$

- Submodular Optimization Problem

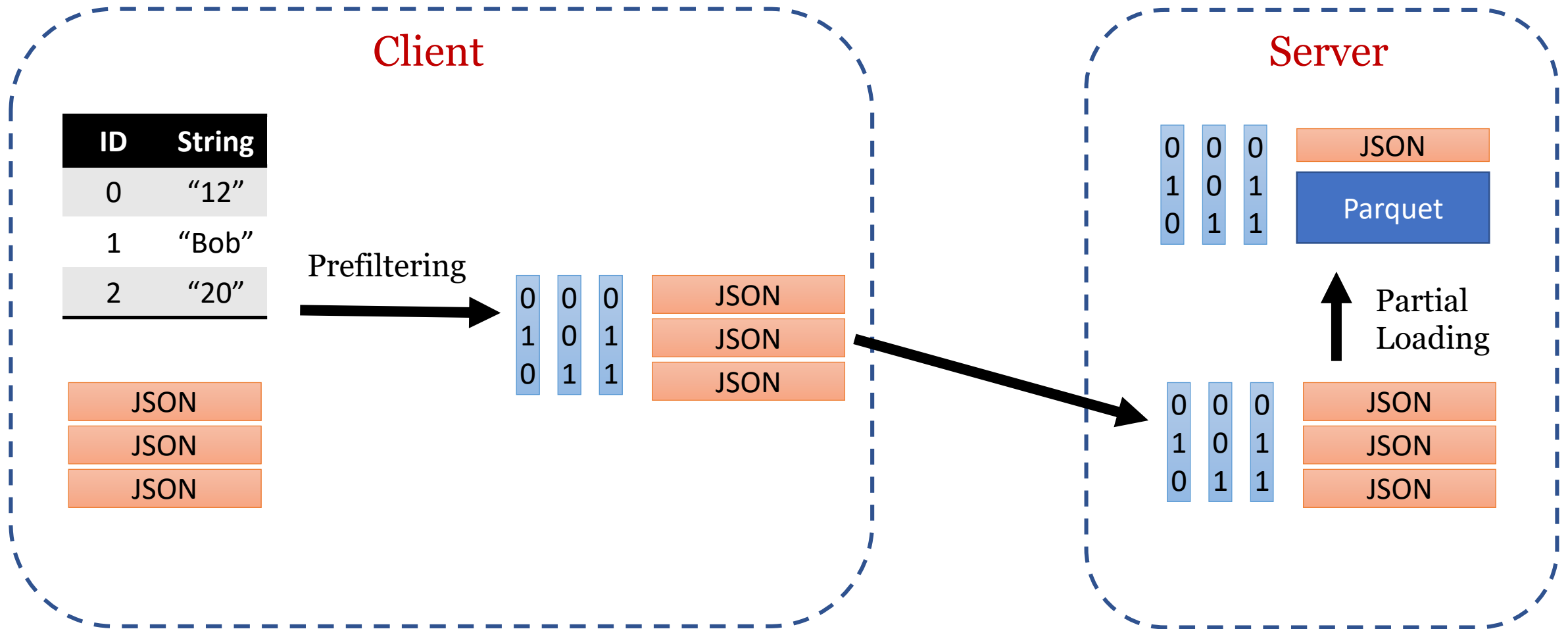
$$\begin{aligned} \max f(S) &= \sum_{q_i \in Q} [1 - \prod_{p_j \in (q_i \cap S)} sel(p_j)] \\ \text{s. t. } \sum_{p_i \in S} cost(p_i) &\leq B \end{aligned}$$

# Overview

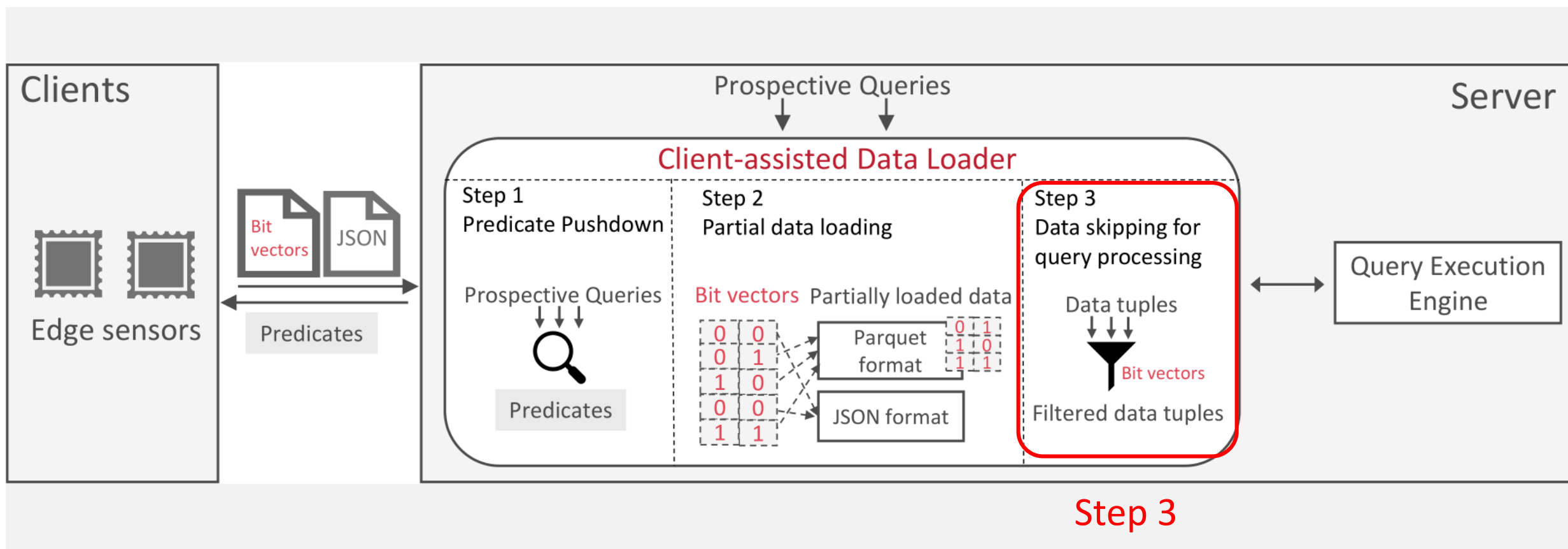




# Partial Data Loading

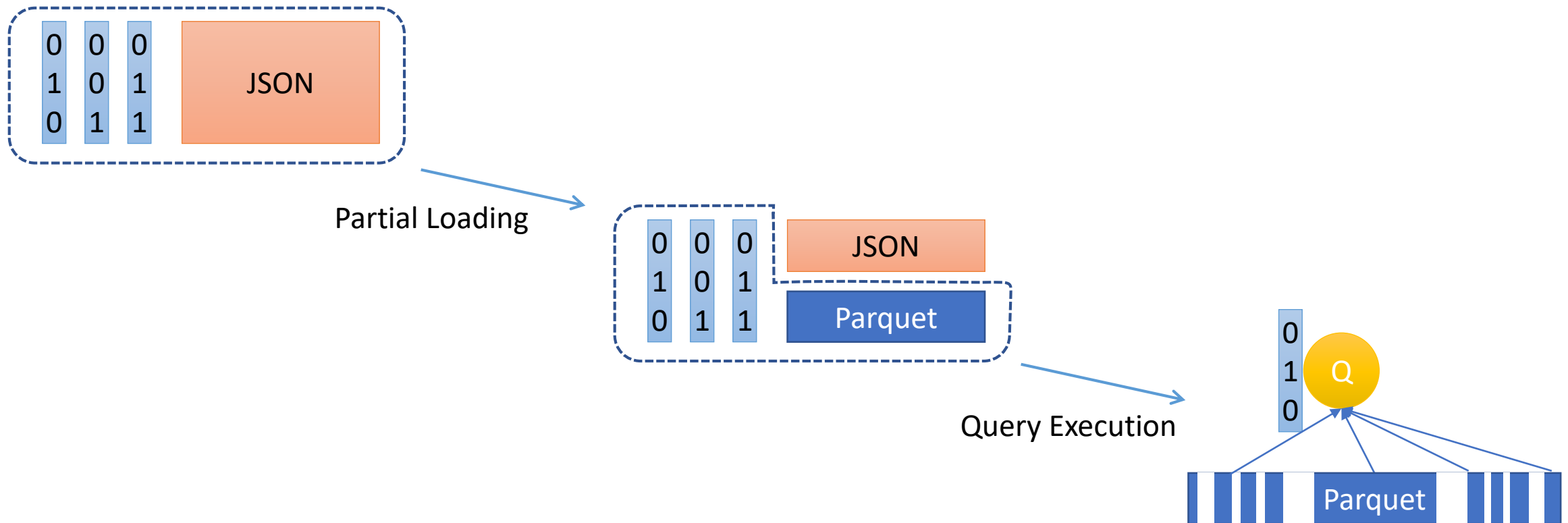


# Overview



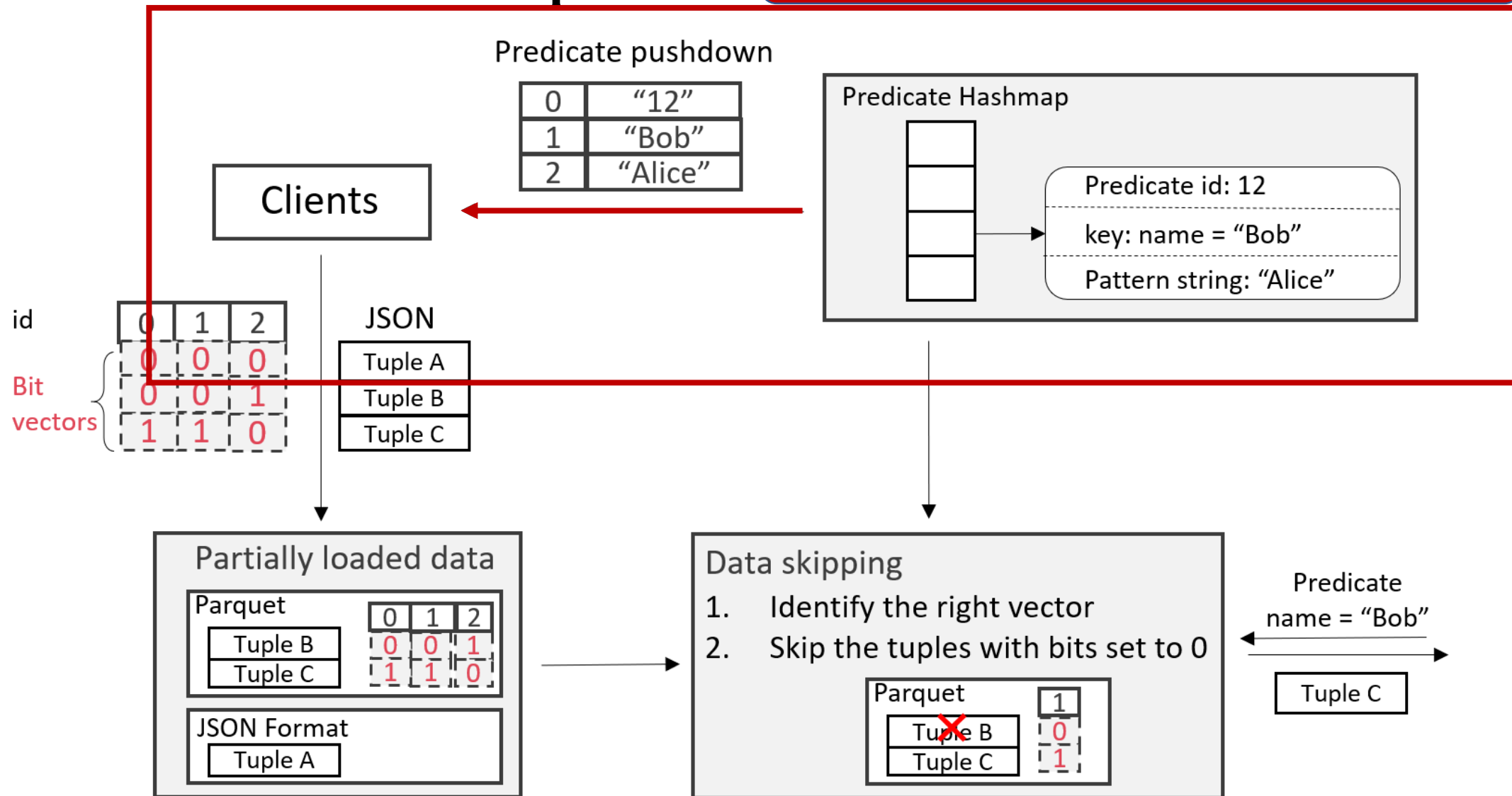
# Data Skipping

- Embedded bit vectors could also help downstream query processing by **skipping irrelevant tuples**



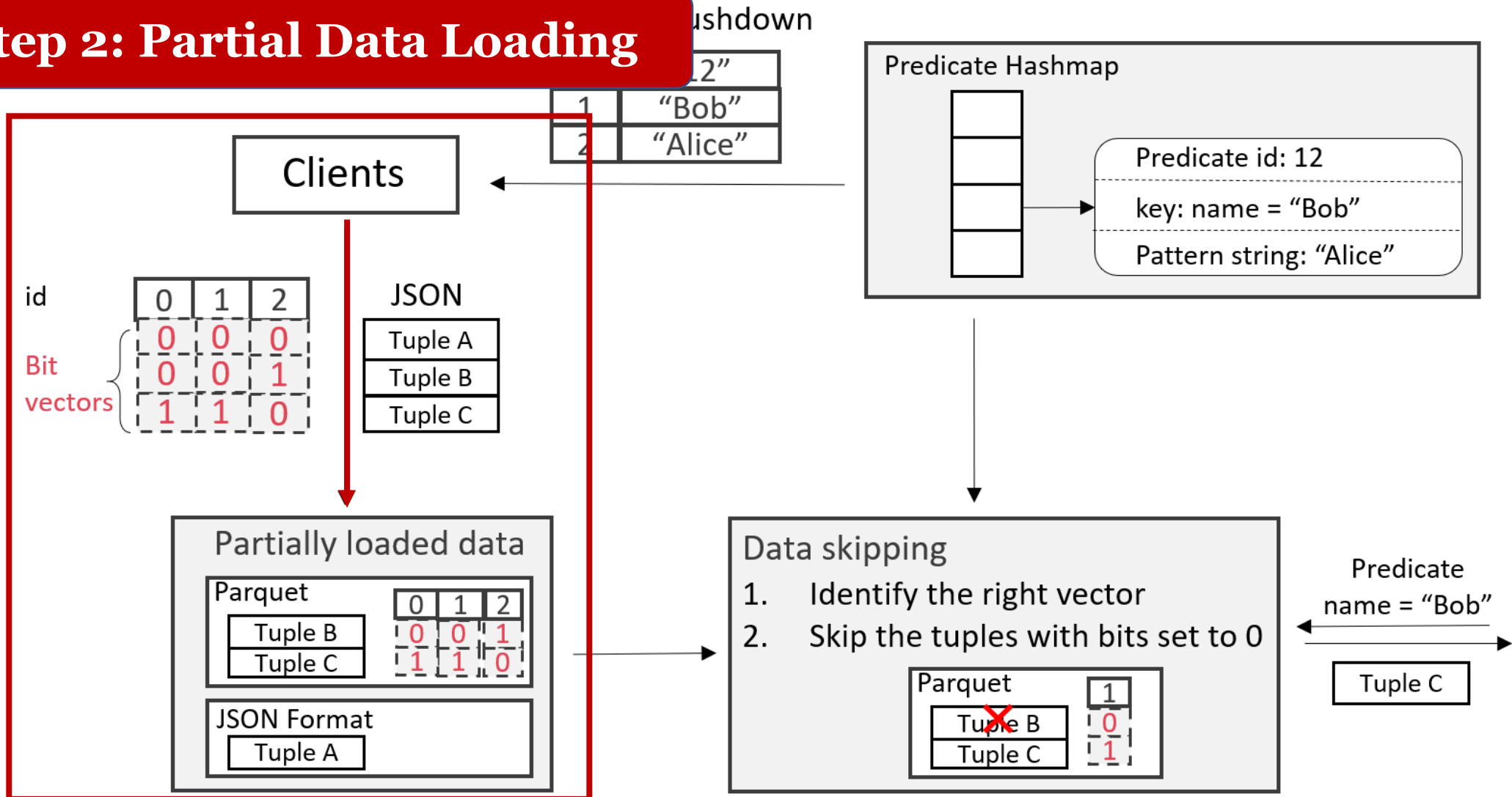
# Workflow Example

## Step 1: Predicate Pushdown

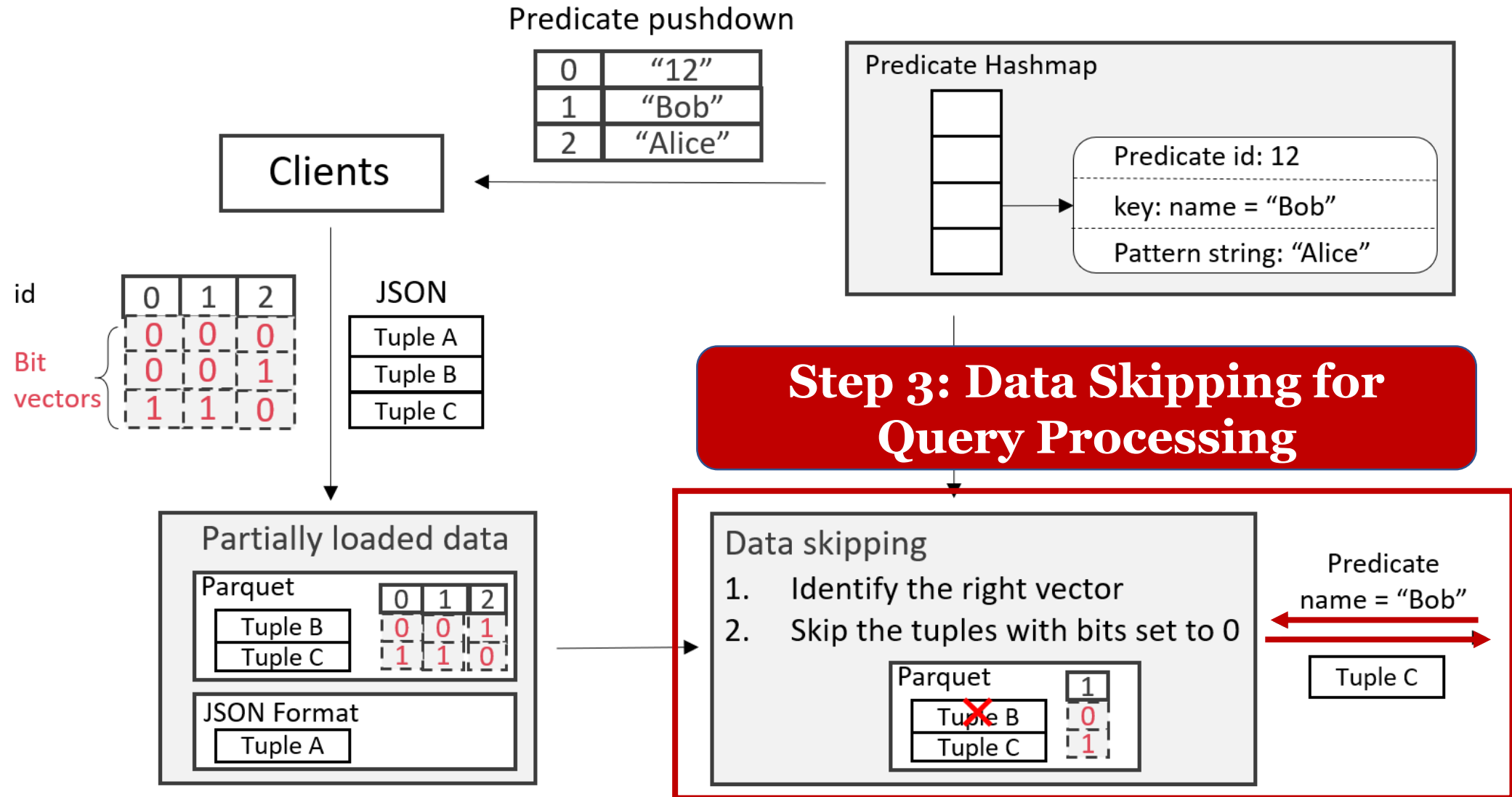


# Workflow Example

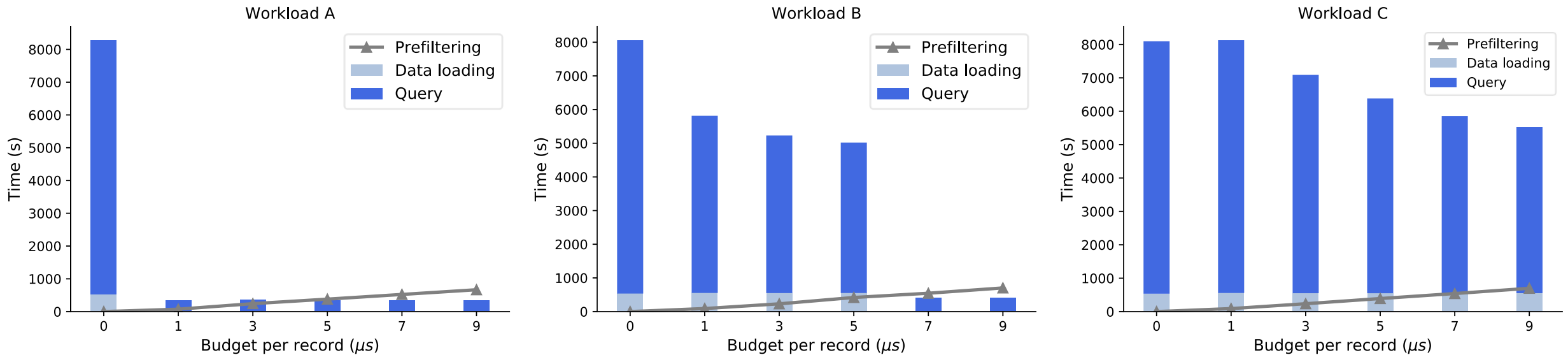
## Step 2: Partial Data Loading



# Workflow Example

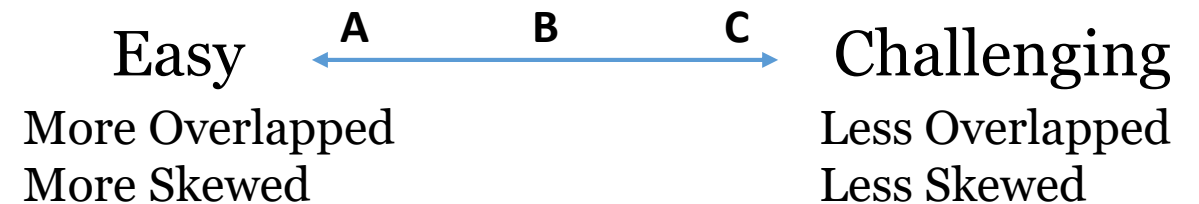


# End-to-End Experiments



- Dataset
  - Windows Event Log
  - 27 GB
  - ~114m Tuples

- Workloads – 200 Queries



# Take-Home Points from *CIAO*

- **Client-Assisted Partial Loading**
  - Clients will help prefilter tuples before loading.
- **Near-Optimal Predicate Selection**
  - Select most beneficial predicate set within limited budget.
- **Data skipping**
  - Downstream query processing will leverage prefiltering results.

Thanks for your attention!

Q&A