# Storing Intermediate Results in Space and Time
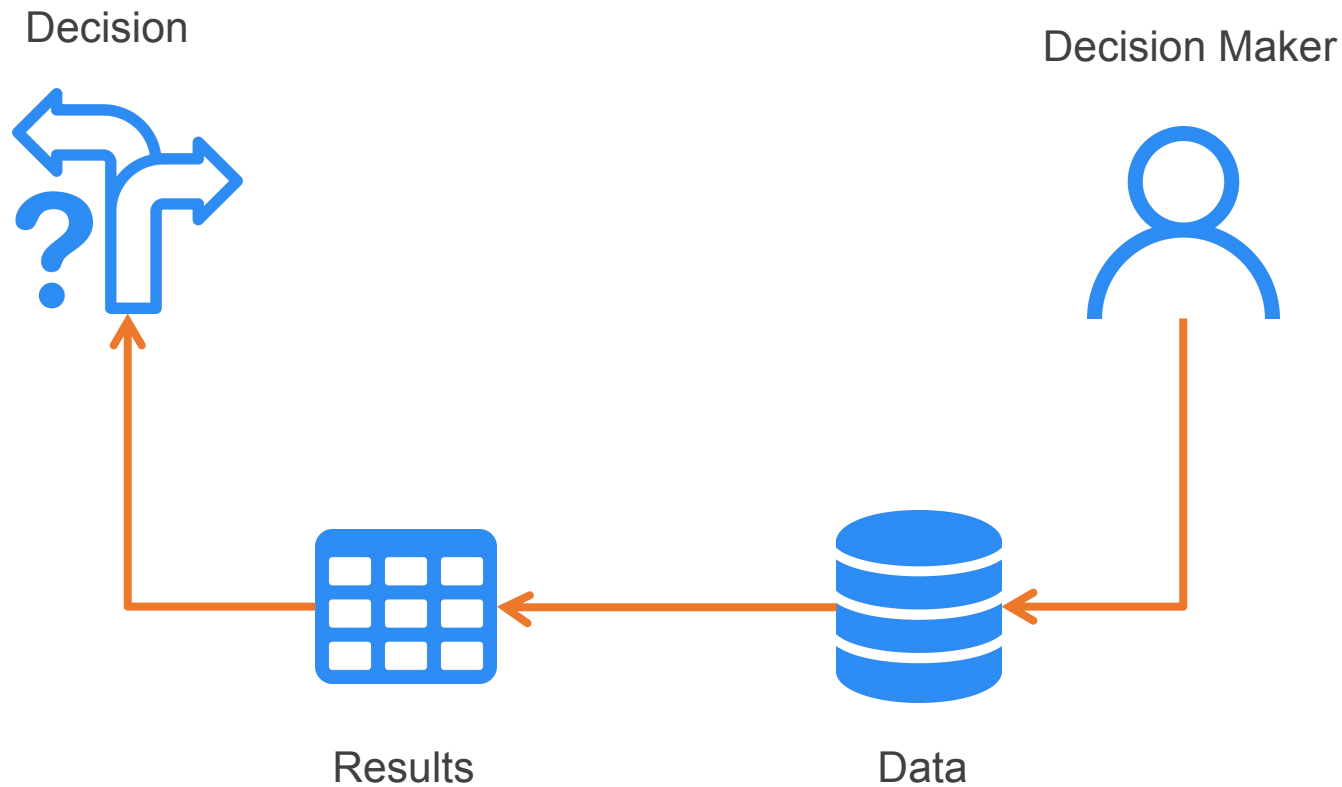
# SQL Graphs in jSQL$_e$

By: *Basem Elazzabi*  (elazzabi@pdx.edu)

Advisors: *David Maier* & *Kristin Tufte*

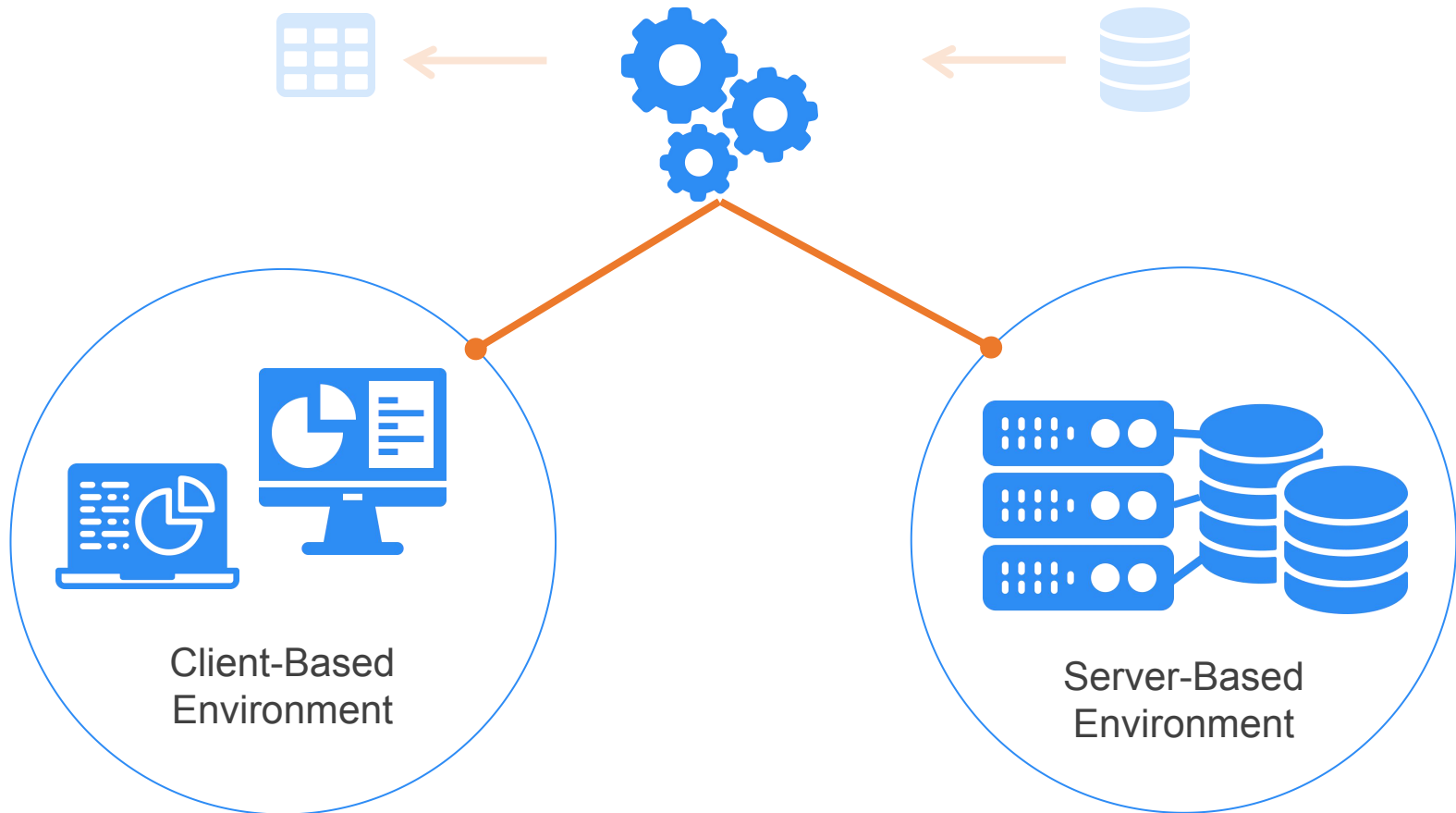Portland State
UNIVERSITY

# The Big Picture

# Data-Driven Decision Making

Decision

Decision Maker

Results

Data

# Data Analysis

Results     Data Analysis     Data

Facilitating the data-analysis process

# Types of data analysis



Client-Based Environment

Server-Based Environment

# Server-Based Environment

User makes
a request

Server receives &
processes the request

Communication
over network

Server produces
the results

User receives
the results

Server sends
back the results

# Client-Based Environment

User imports a data set into a PC or laptop

Data Set

User uses multiple tools to analyze the data

Final Results

Data Analysis Tools

# Client-Based Environment



User imports a data set into a PC or laptop

Data Set

User uses multiple tools to analyze the data

Final Results

Data Analysis Tools

# The Problem

# Typical Data-Analysis Process



Import data

Perform data manipulations

Data Set

Excel / Spreadsheets

Create some charts

# Typical Data-Analysis Process



Data Set

Import data

Create visualizations

Export visualizations

PNG     JPEG

# Typical Data-Analysis Process
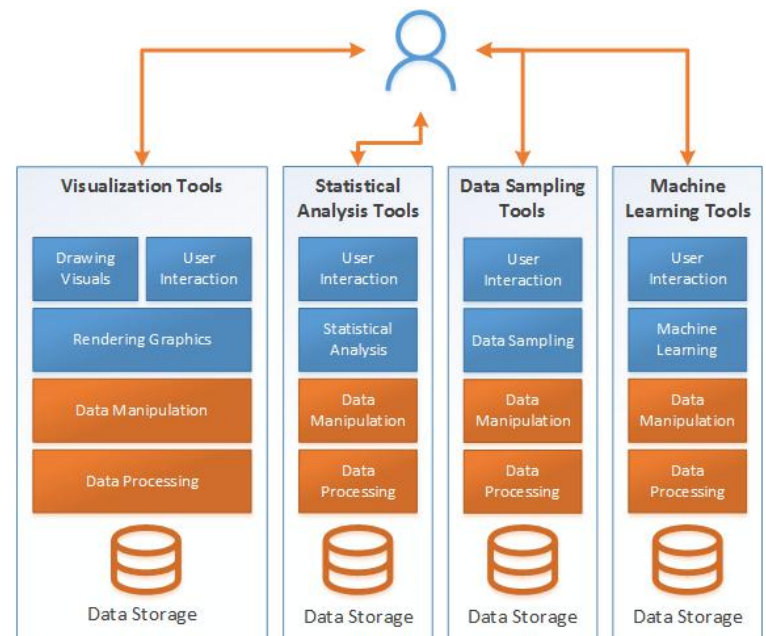


Manual data movement & conversion
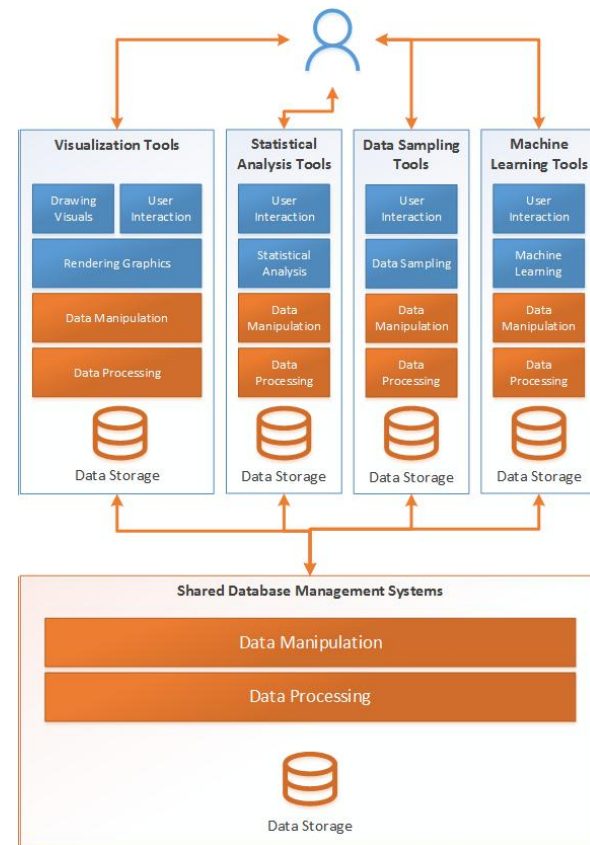
# Manual Data Movement

# Manual Data Movement

- Wastes a lot of time

- Takes a lot of effort

- Requires technical skills

- Wastes space

- Redoing computations
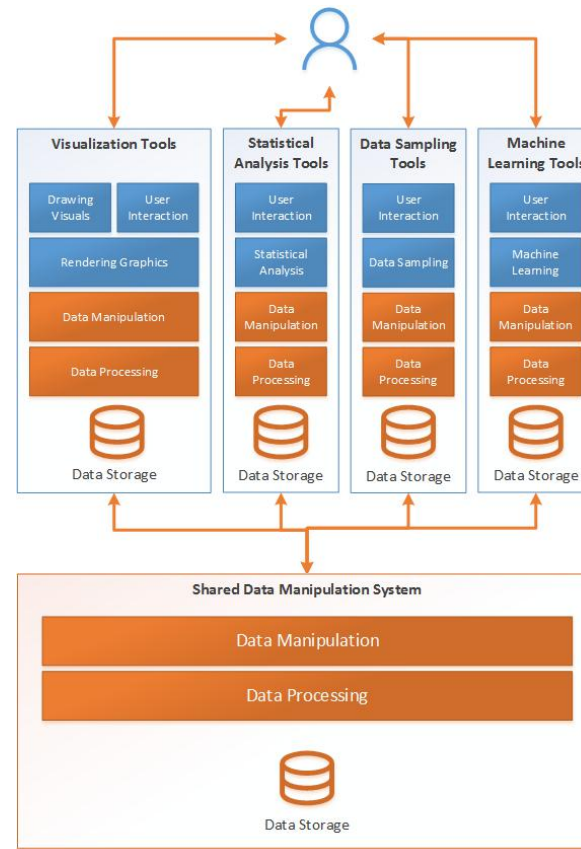
- Difficult to inspect results

# Shared Data-Manipulation Systems

- Still wastes a lot of time

- Still takes a lot of effort

- Still requires technical skills

- Still wastes space

- Still redoing computations
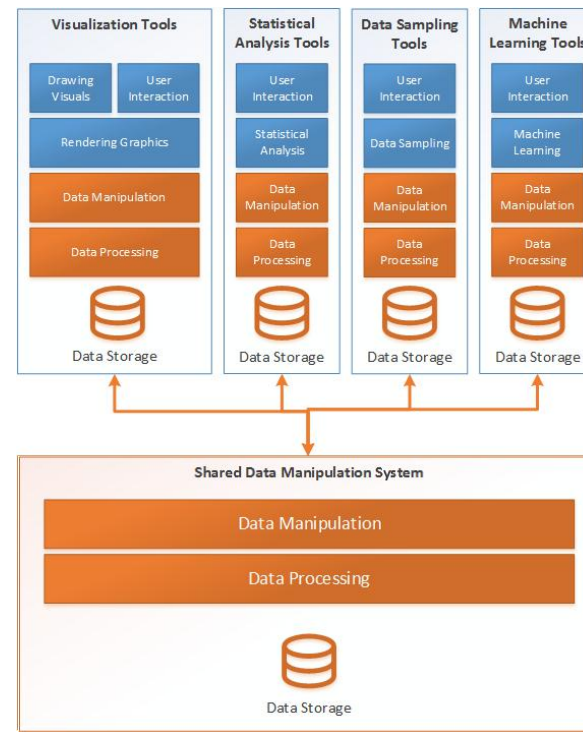
- Still difficult to inspect results

# The Ideal Environment
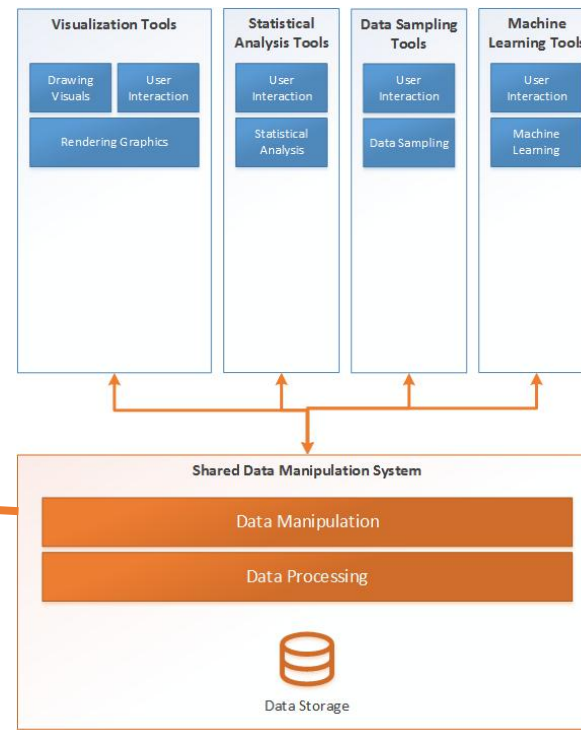
# The Ideal Environment

- Eliminate data movement and conversion
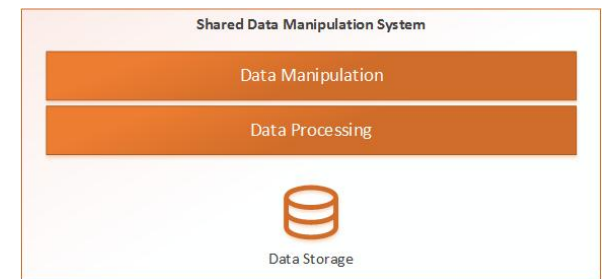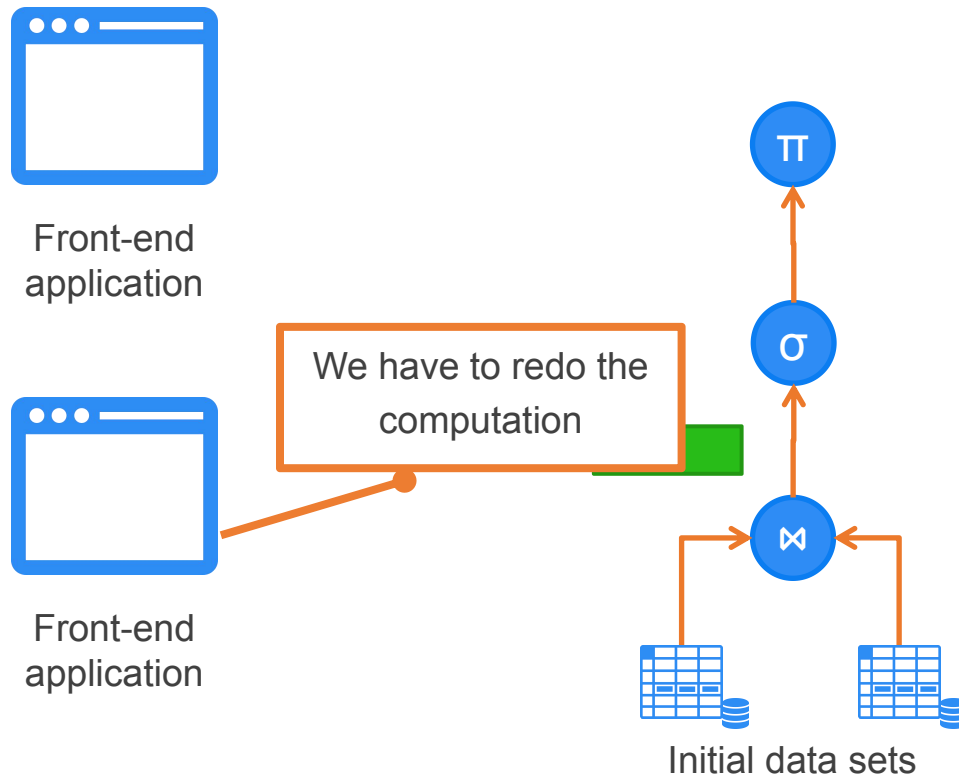
# The Ideal Environment

- Eliminate data movement and conversion
- Factor out the data manipulation process and storage into the shared system.
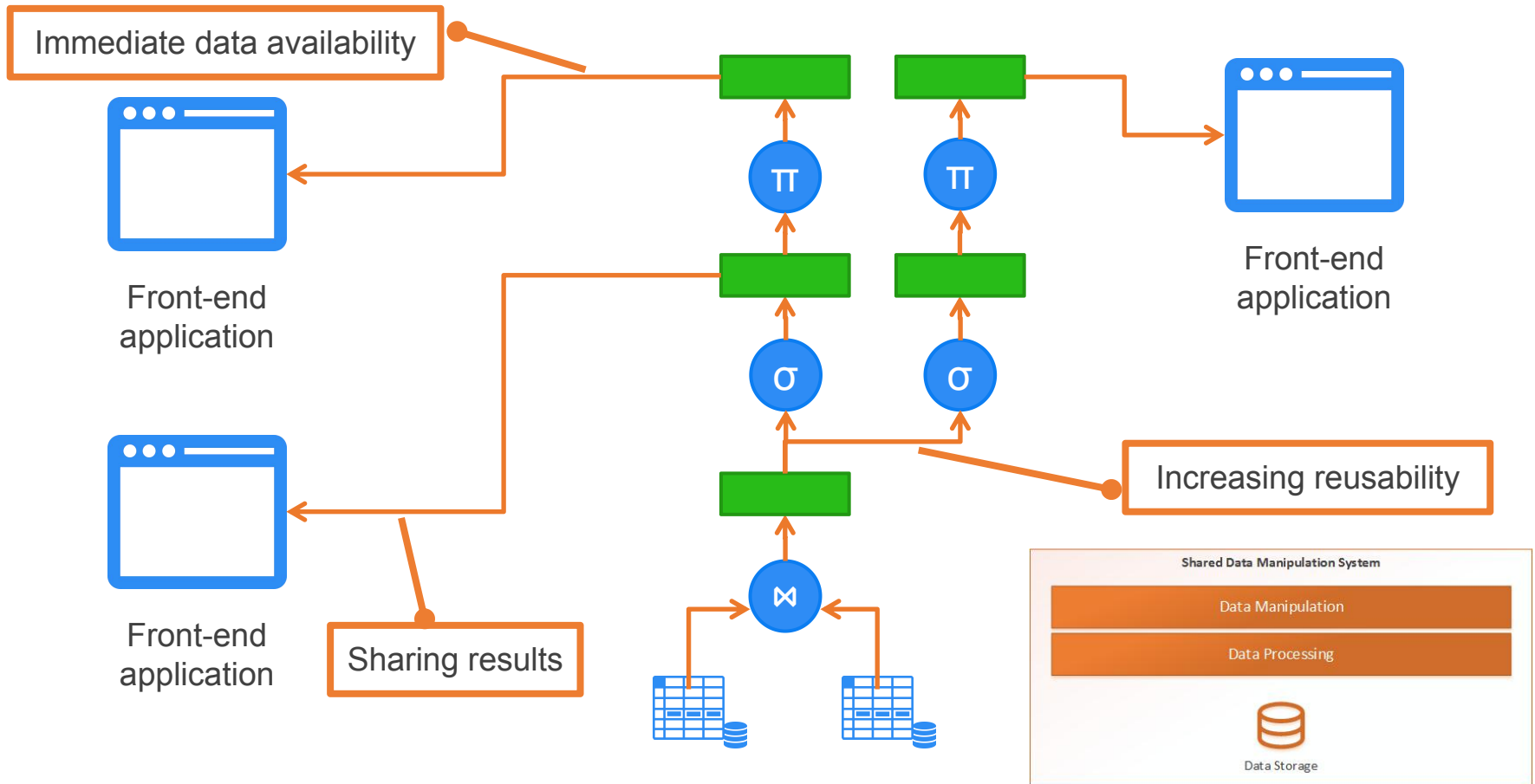
How can we build this system?

# The Data Model

# The Current Data Model



Front-end application

Front-end application

We have to redo the computation

π

σ

⋈

Initial data sets

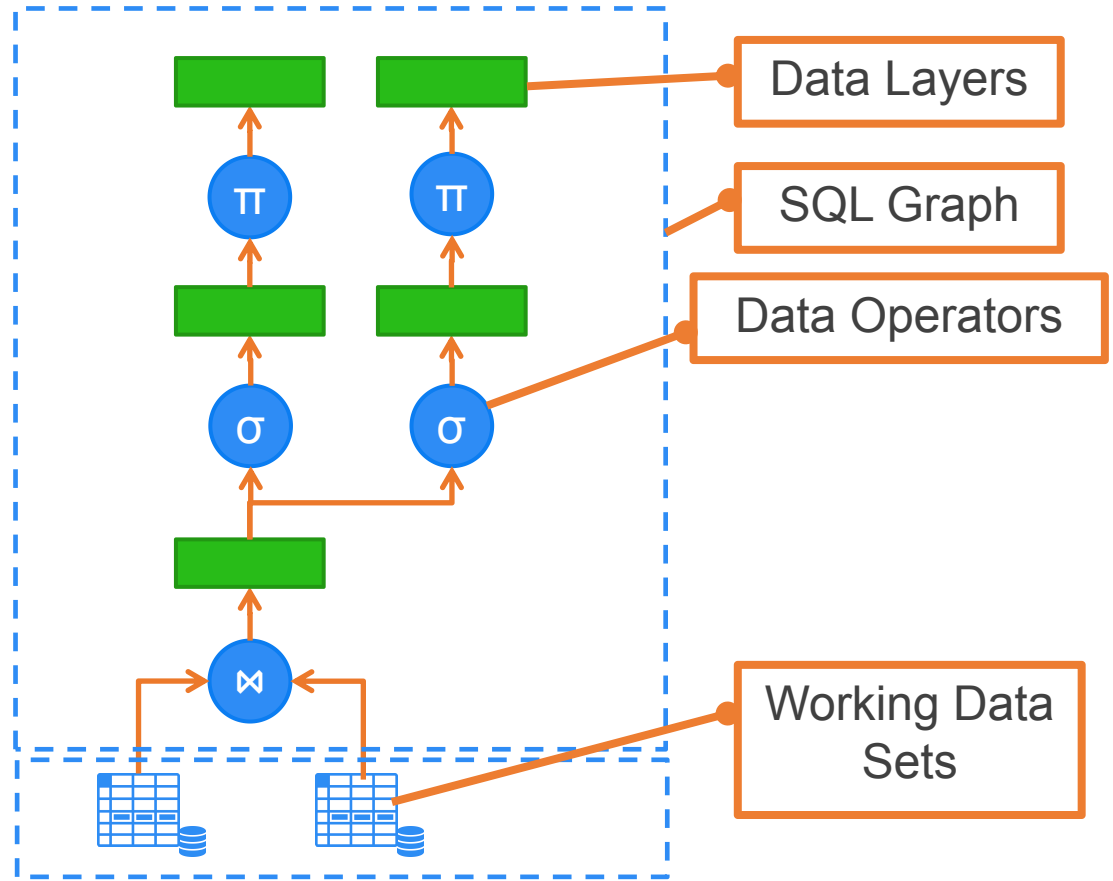**Shared Data Manipulation System**

Data Manipulation

Data Processing
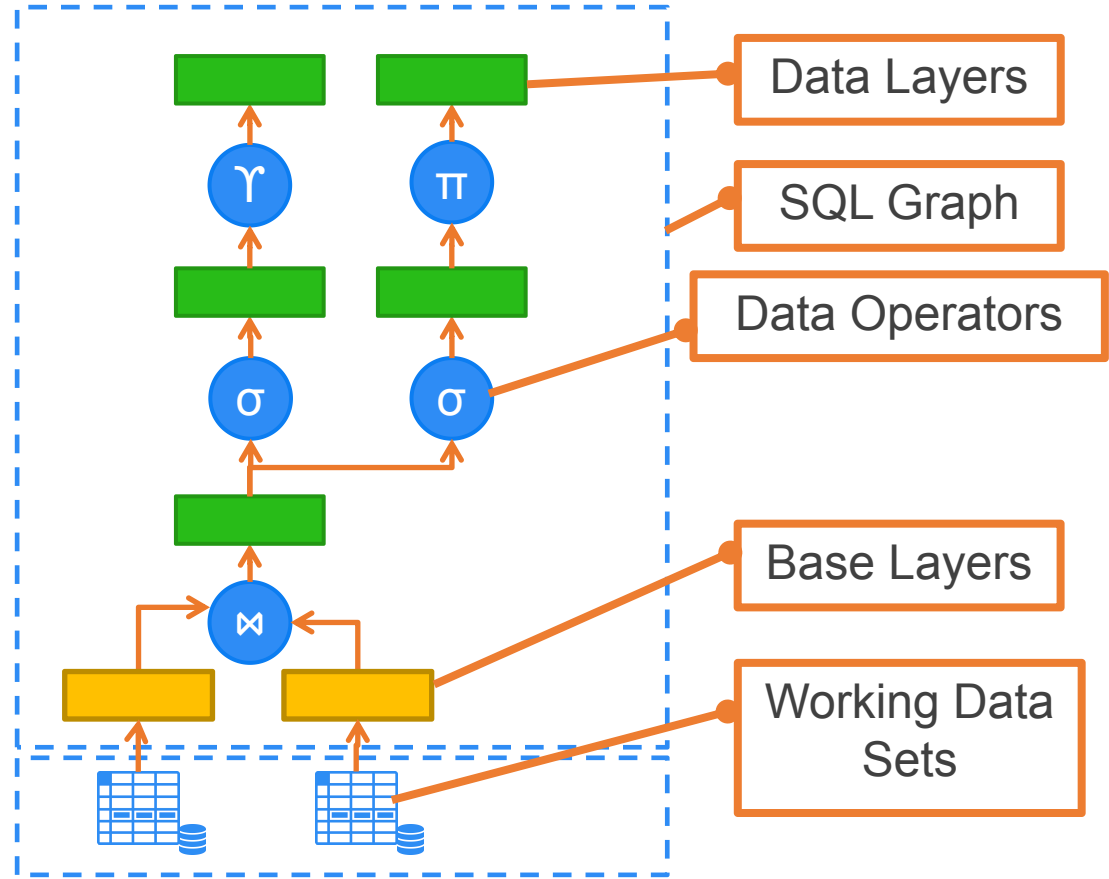
Data Storage

# The Proposed Data Model

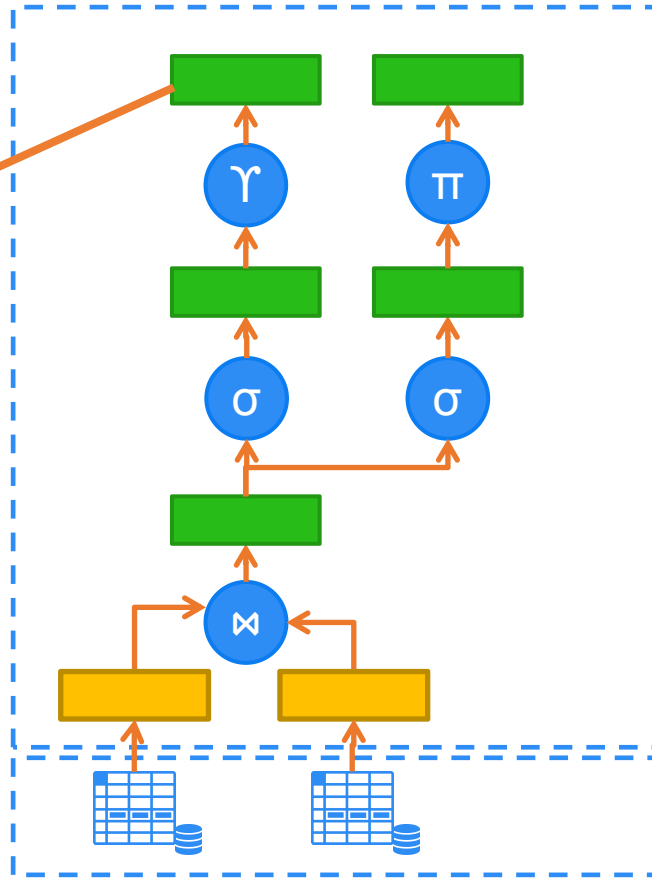# The Proposed Data Model

# The Proposed Data Model

- Immediate data availability
- Share results
- Eliminate data movement
- Eliminate data conversion
- Increase reusability

# The Challenges

# The Challenges

1. How to store intermediate results efficiently.

2. How to provide data accessibility with interactive speed.
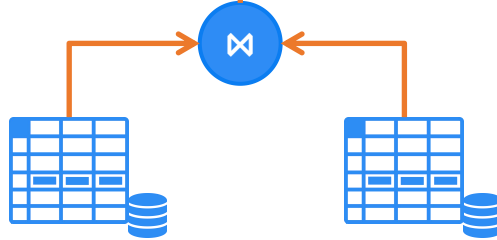
# Storage Cost - Materialization

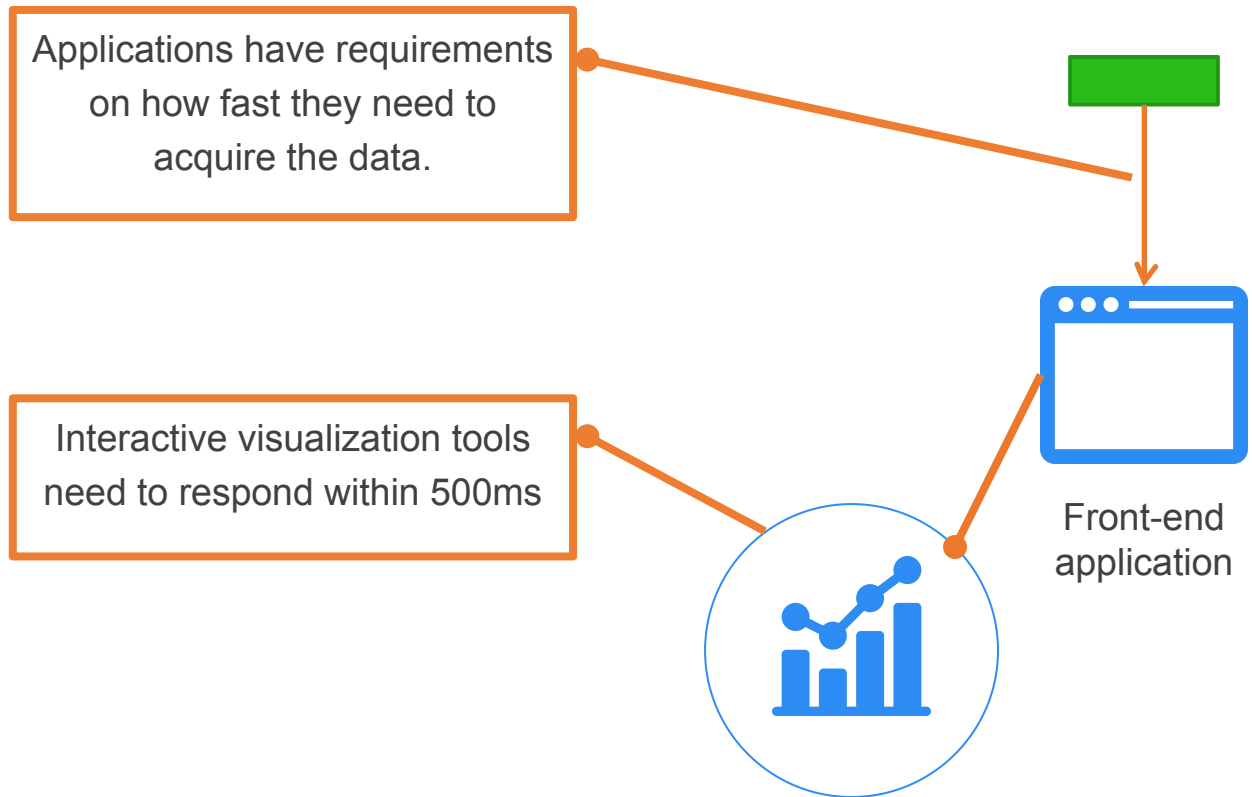Data materialization is very inefficient

> 8 GB

> 4 GB

4 GB          100 KB

# Data-Access Time

Applications have requirements on how fast they need to acquire the data.

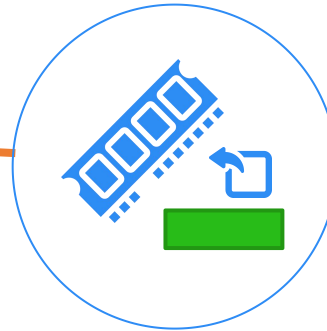Interactive visualization tools need to respond within 500ms
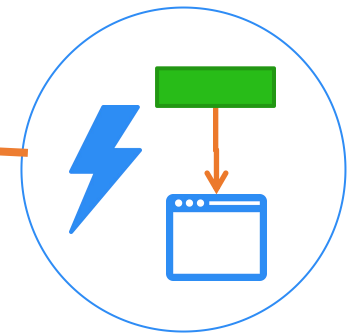
Front-end application

# Research Goals

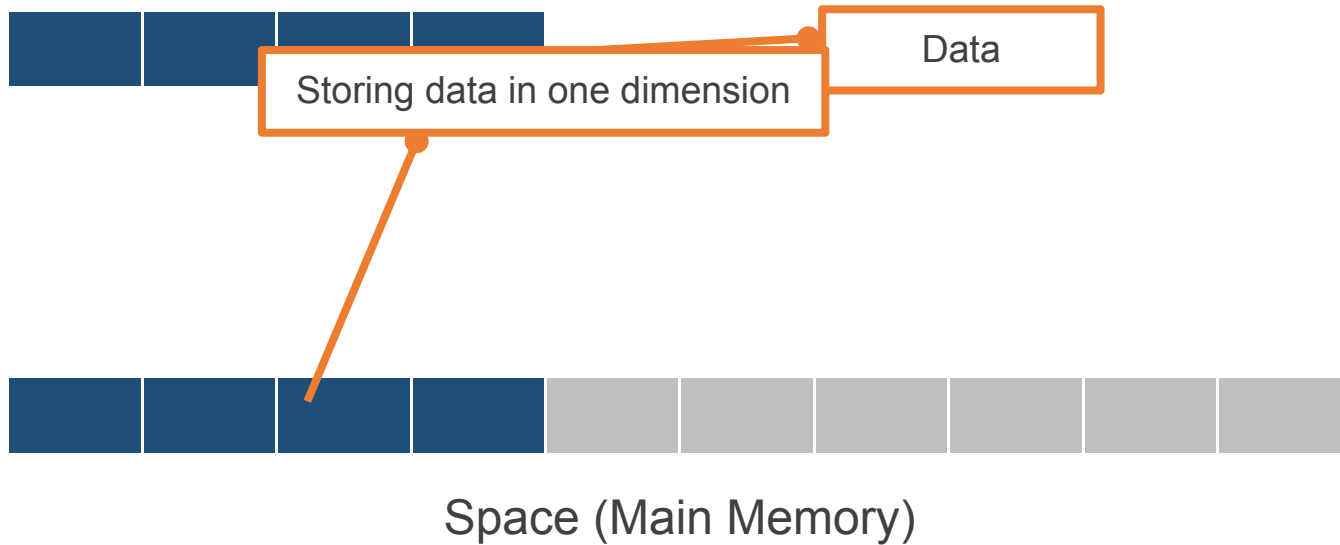Store the data and all or most intermediate results in main memory

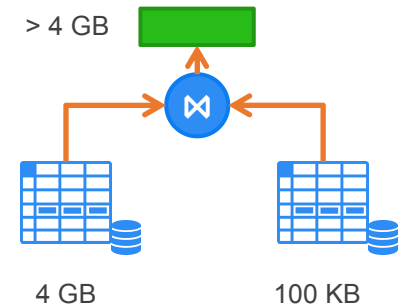Provide data accessibility with interactive speed

# Storing Data in Space and Time

# Storing Data in One Dimension

Data

Storing data in one dimension

Space (Main Memory)

# Storing Data in Two Dimensions



Time (CPU)

Space (Main Memory)

8GB

> 4 GB

4 GB

100 KB

# Storing Data in Two Dimensions



Interactive-speed threshold

Time (CPU)

Space (Main Memory)

8GB

> 1.6 GB

1.6 GB

40 KB

# Storing Data in Two Dimensions

Interactive-speed threshold

Time (CPU)

Space (Main Memory)

8GB

# The Solution

# Space Cost

Using materialization, total space cost of intermediate results is > 4GB

Time (CPU)

Space (Main Memory)

8GB

> 4 GB

4 GB

100 KB

# Space Cost

Using block referencing, total space cost of intermediate results is 240MB

240 MB

4 GB    100 KB

Time (CPU)

Space (Main Memory)

8GB

# Space Cost of Working Data Sets



The space cost of the working data sets is still about 4GB

Leaves only 1GB for data analysis

Time (CPU)

Space (Main Memory)

240 MB

4 GB

100 KB

8GB

# The New Storage Engine

Byte Array

Bookkeeping

# Space Cost of Working Data Sets



Time (CPU)

Space (Main Memory)

8GB

240 MB

4 GB

100 KB

# Space Cost of Working Data Sets

The new space cost of the working data sets is about **900 MB** instead of **4 GB**

Time (CPU)

Space (Main Memory)

8GB

240 MB

900 MB          100 KB

# Experiment

# Realistic Use-Case

# Experiment Goals

- Perform a realistic data-analysis use-case using **jSQL$_e$**.

- Keep the results of all intermediate results and the original

  data set in main memory.

- Simulate the same analysis in three other well-known

  systems and compare the space and time costs.

# Data-Analysis Objective



For a given **stop Id**, a given **bus/train route**, and a given **schedule time**, what will be the **actual arrival time** for that bus/train?

Data Analysis

Historical Public-Transportation Data

Model

12:45 pm

Prediction

Accuracy within **± 3 minutes** from actual arrival time

# Data-Analysis Objective



Data Analysis

Model

Historical Public-Transportation Data

12:45 pm

Prediction

# Data-Analysis Objective

# Experiment Setup

- We used **6 months** of historical data from **TriMet.**

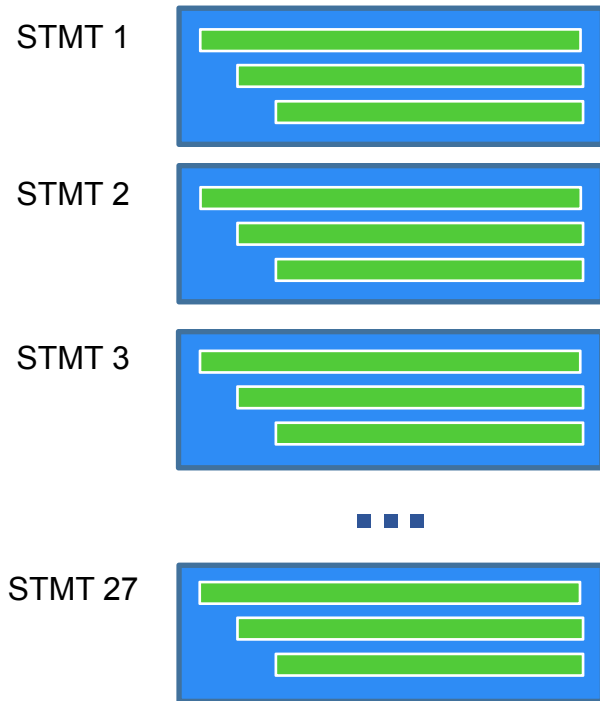- We used a desktop computer

  - RAM: **8GB**

  - CPU: **4 cores, i5, 3.5GH**

- Space limit for data analysis is **6GB**

- We ran all four systems on a single core (one thread).

# Experiment Setup

**SQL (**declarative**)**

**jSQL (**imperative**)**

STMT 1

STMT 2

STMT 3

STMT 27

Step 1
Step 2
Step 3

Step 178

# Experiment Setup

**SQL (**declarative**)**
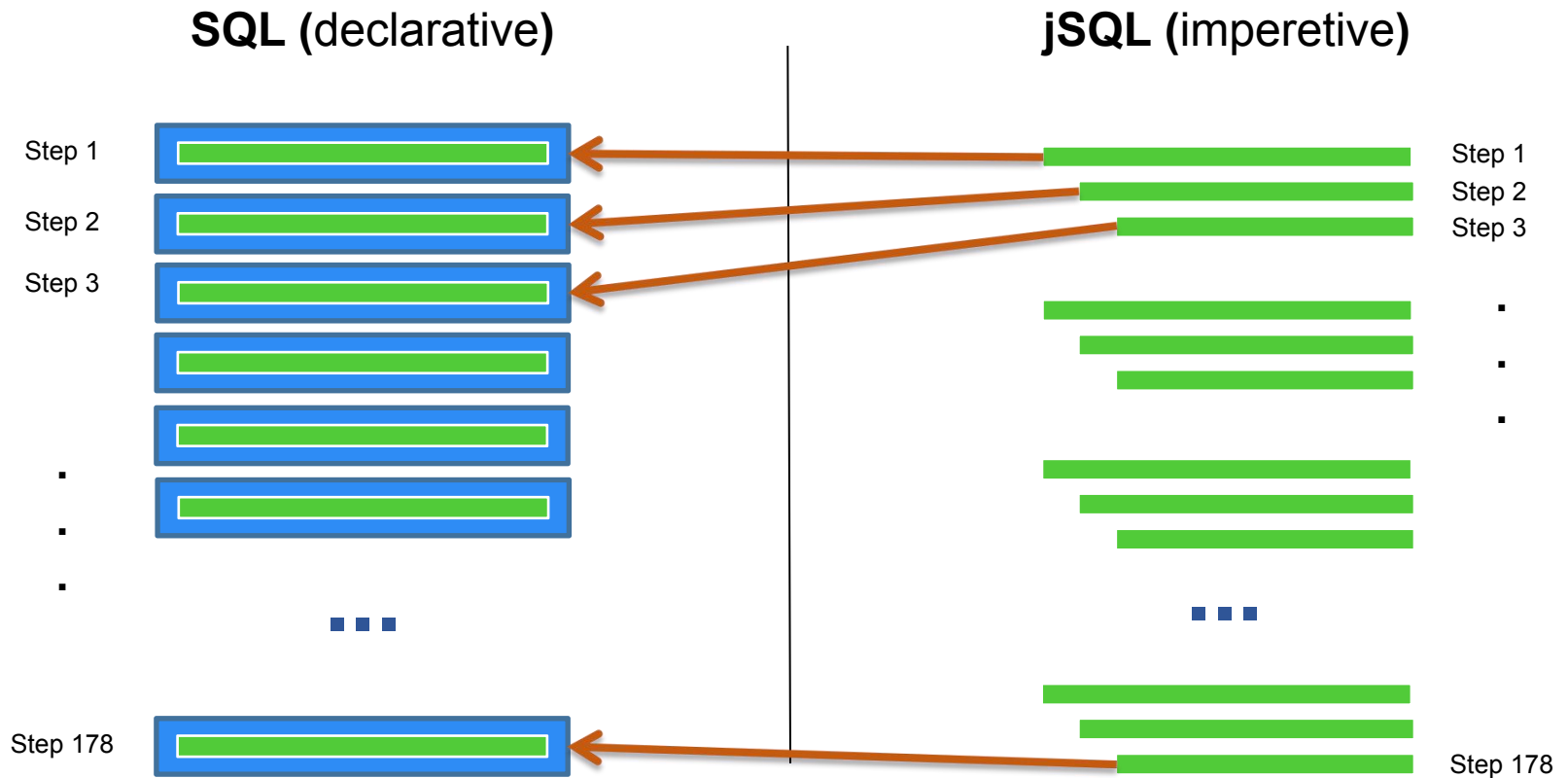
```
SELECT
    service_date, route_number,
    schedule_time, arrive_time
FROM
    historical_data
WHERE
    service_date = '2021-01-29';
```

**jSQL (**imperative**)**

```
L1 = SELECT historical_data
    WHERE
        service_date = '2021-01-29';


L2 = PROJECT L1 WITH
    service_date, route_number,
    schedule_time, arrive_time;
```

# Experiment Setup

**SQL (**declarative**)**

**jSQL (**imperetive**)**

Step 1

Step 2

Step 3

Step 178

Step 1

Step 2

Step 3

Step 178

# Experiment Setup

**PostgreSQL (**declarative**)**

```
CREATE TABLE L1 AS
SELECT * FROM historical_data
WHERE
   service_date = '2021-01-29';
```
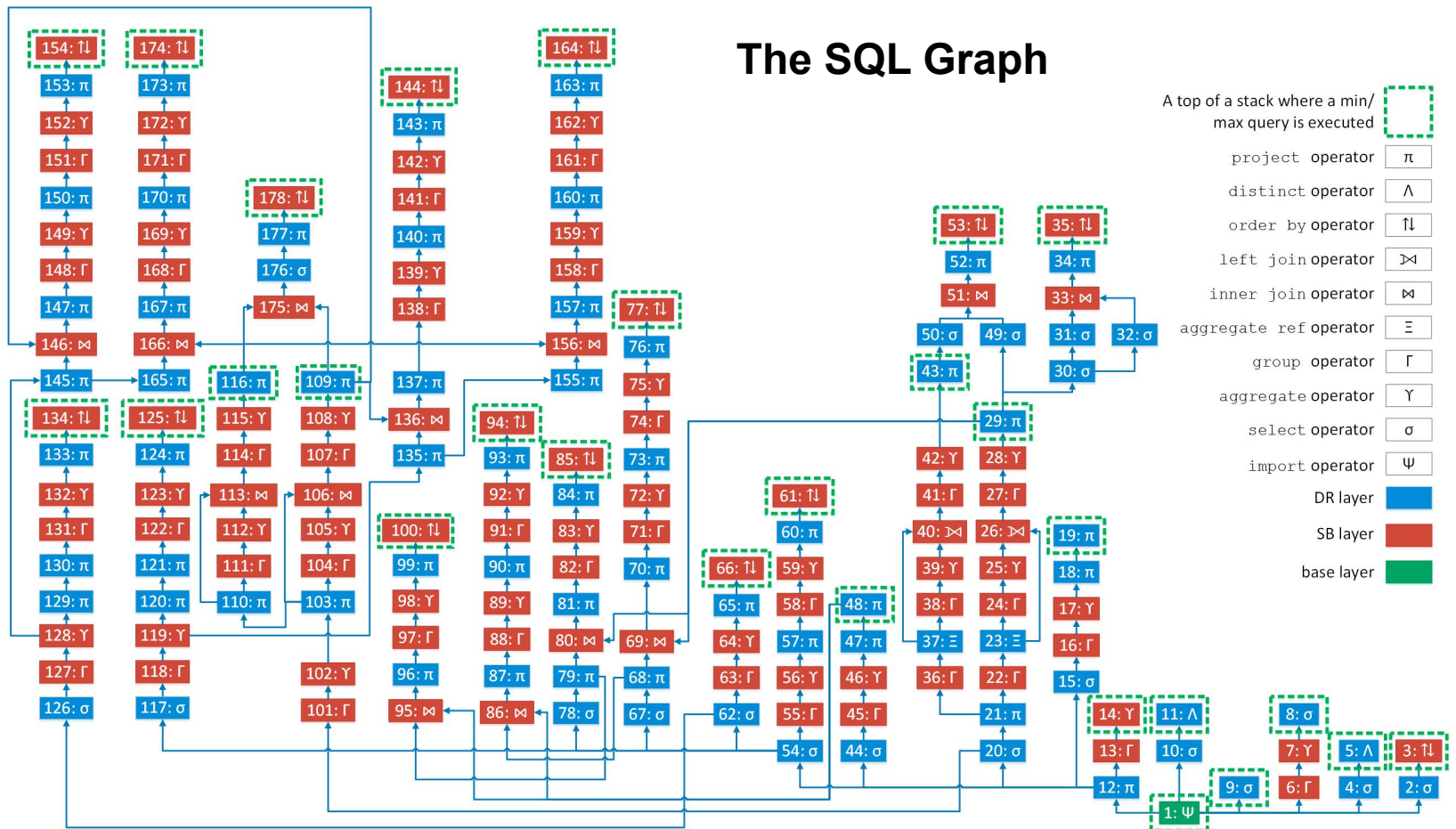
```
CREATE TABLE L2 AS
SELECT
   service_date, route_number,
   schedule_time, arrive_time
FROM L1;
```

**jSQL (**imperative**)**

```
L1 = SELECT historical_data
   WHERE
      service_date = '2021-01-29';
```

```
L2 = PROJECT L1 WITH
   service_date, route_number,
   schedule_time, arrive_time;
```
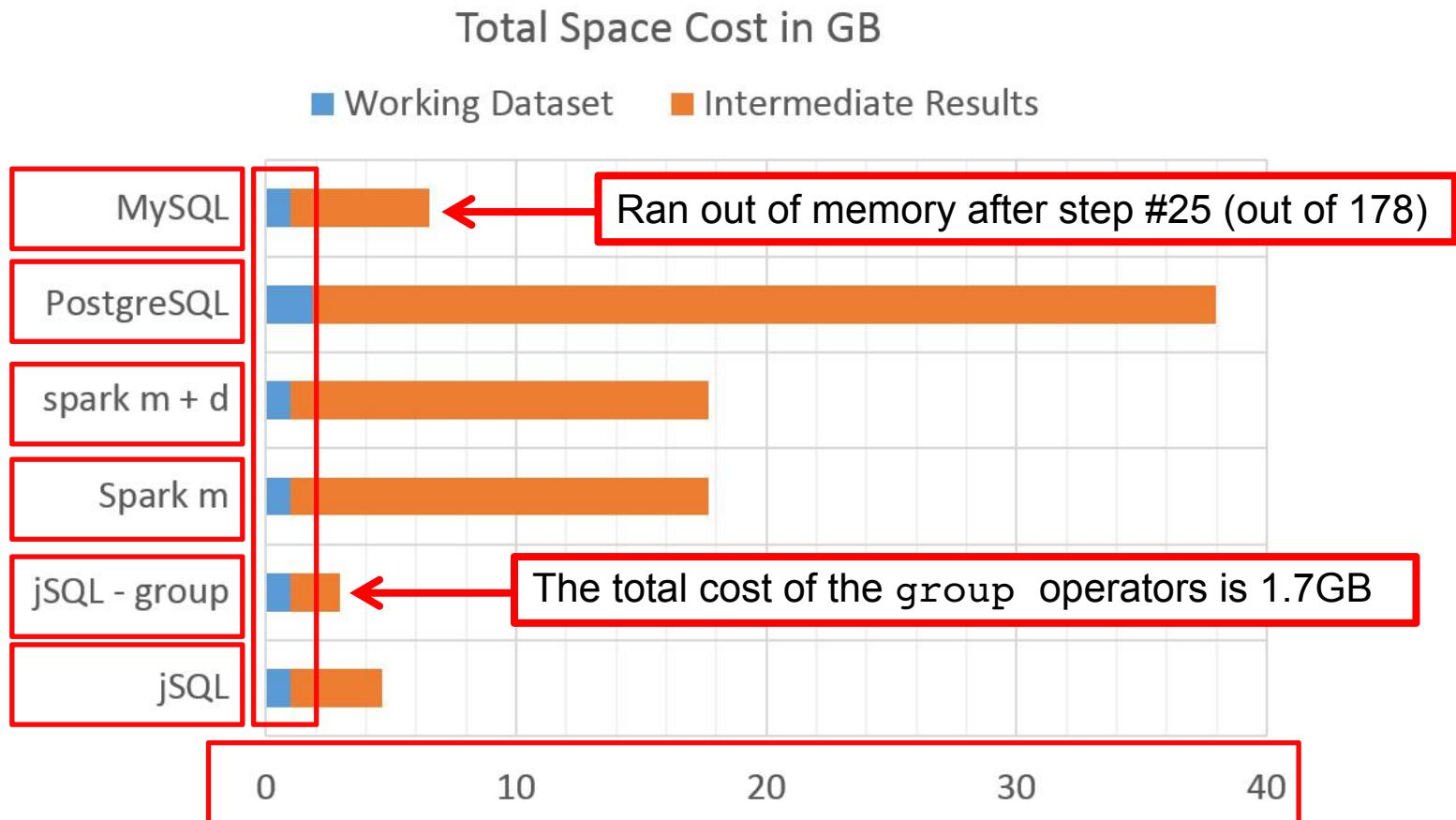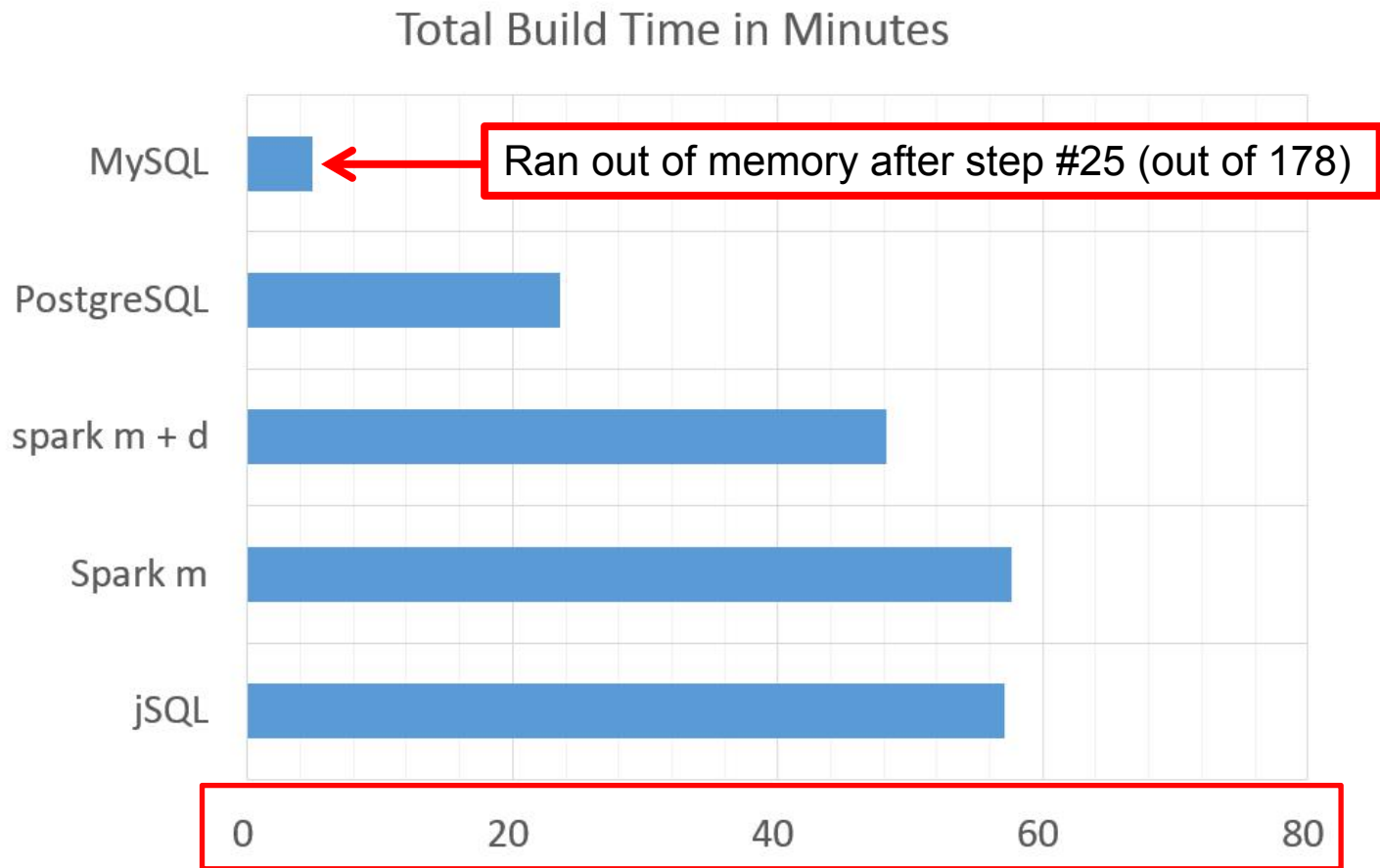
# Experiment Setup



**The SQL Graph**

# Experiment Results

# Realistic Use-Case

# Total Space Cost

## Total Space Cost in GB

■ Working Dataset    ■ Intermediate Results

Ran out of memory after step #25 (out of 178)

The total cost of the `group` operators is 1.7GB

MySQL, PostgreSQL, spark m + d, Spark m, jSQL - group, jSQL

0    10    20    30    40

# Total Build Time

## Total Build Time in Minutes



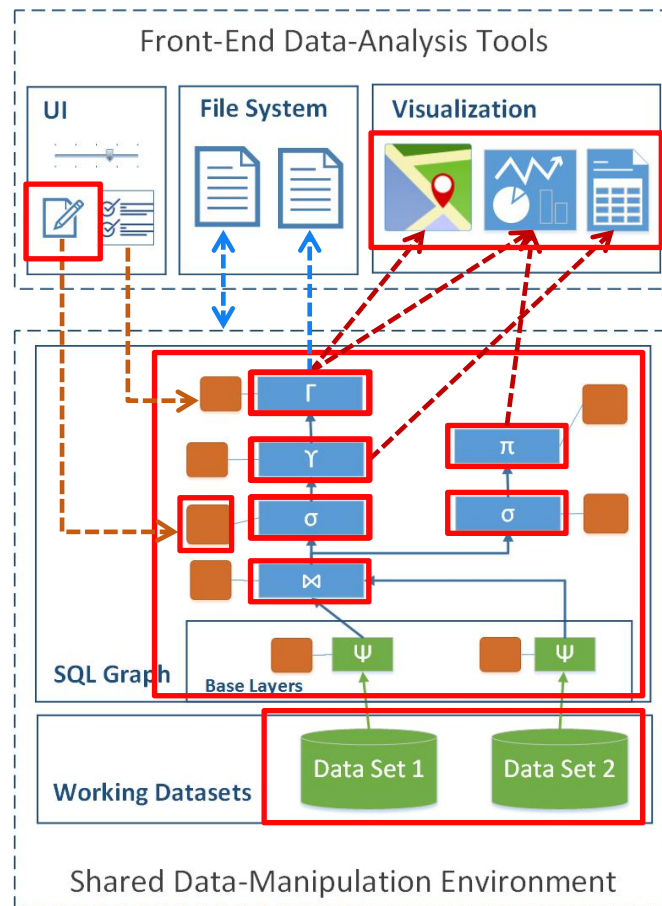Ran out of memory after step #25 (out of 178)

# Lessons Learned

- $jSQL_e$ achived 92% reduction in space cost compared to PostgreSQL, while spending twice the build time.

- $jSQL_e$ achived 83% reduction in space cost compared to Spark, while spending more or less the same build time .
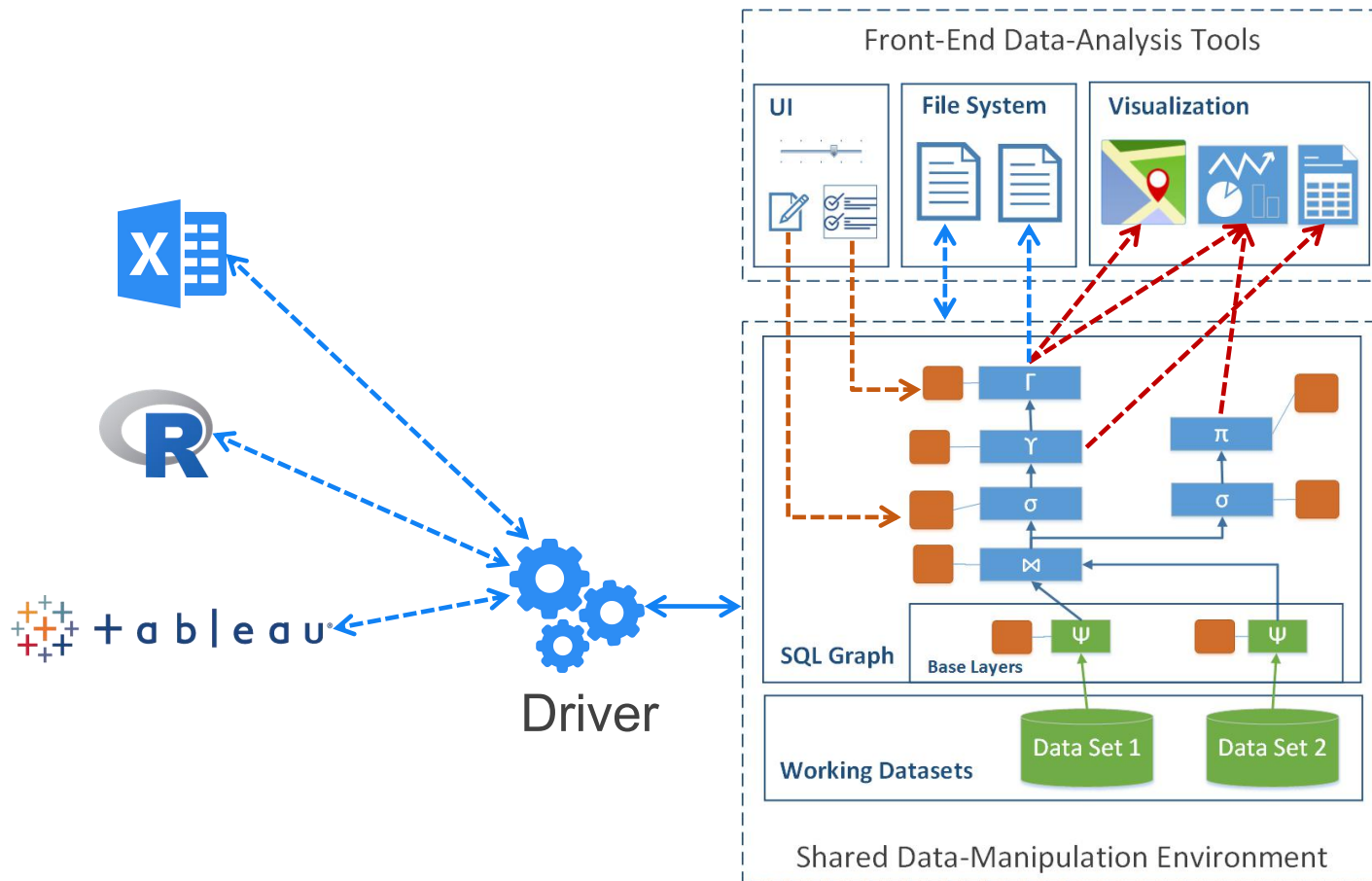
- MySQL? Just don't use it.

# Lessons Learned

- Using current systems to do exploratory data analysis is
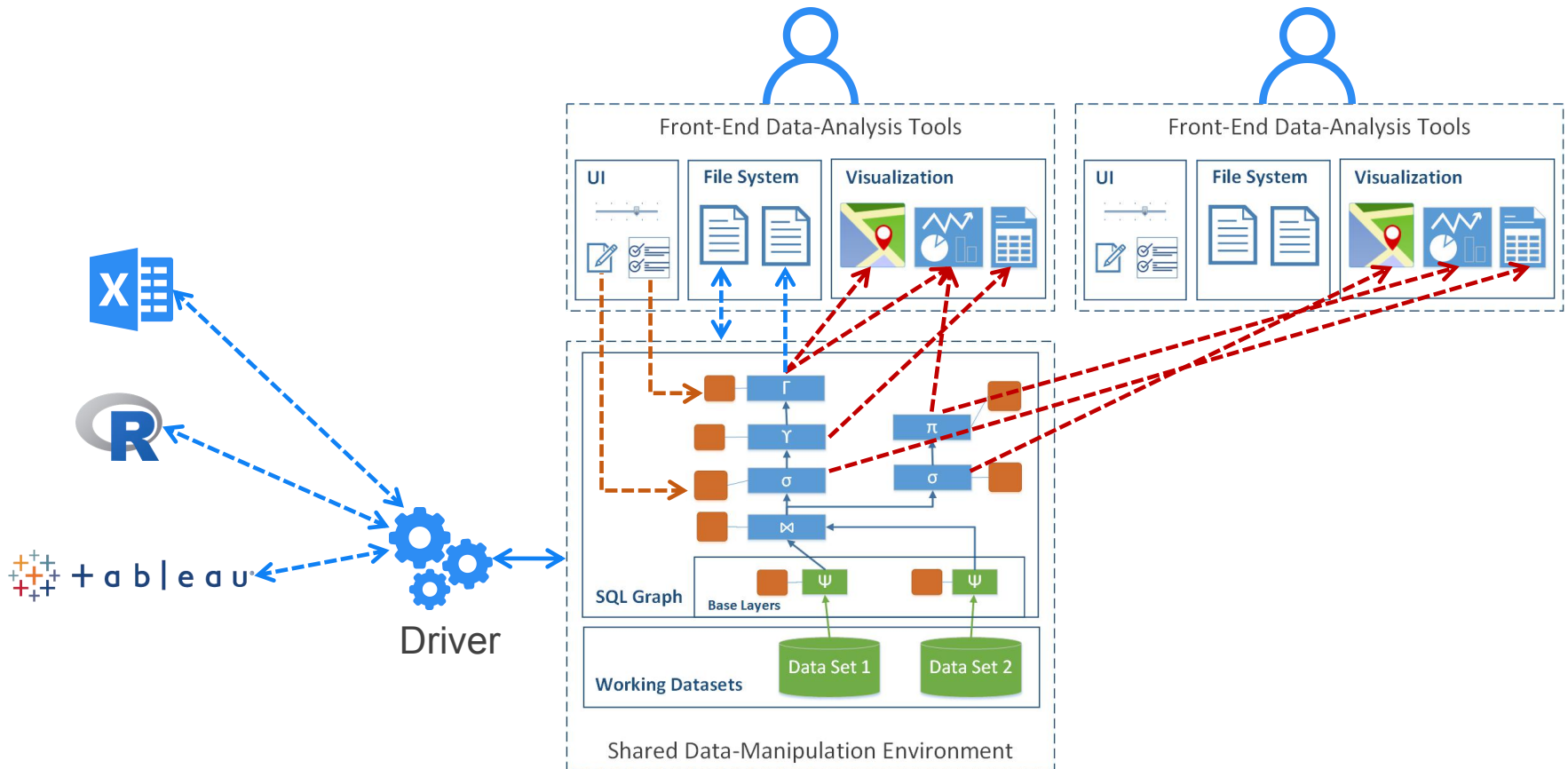
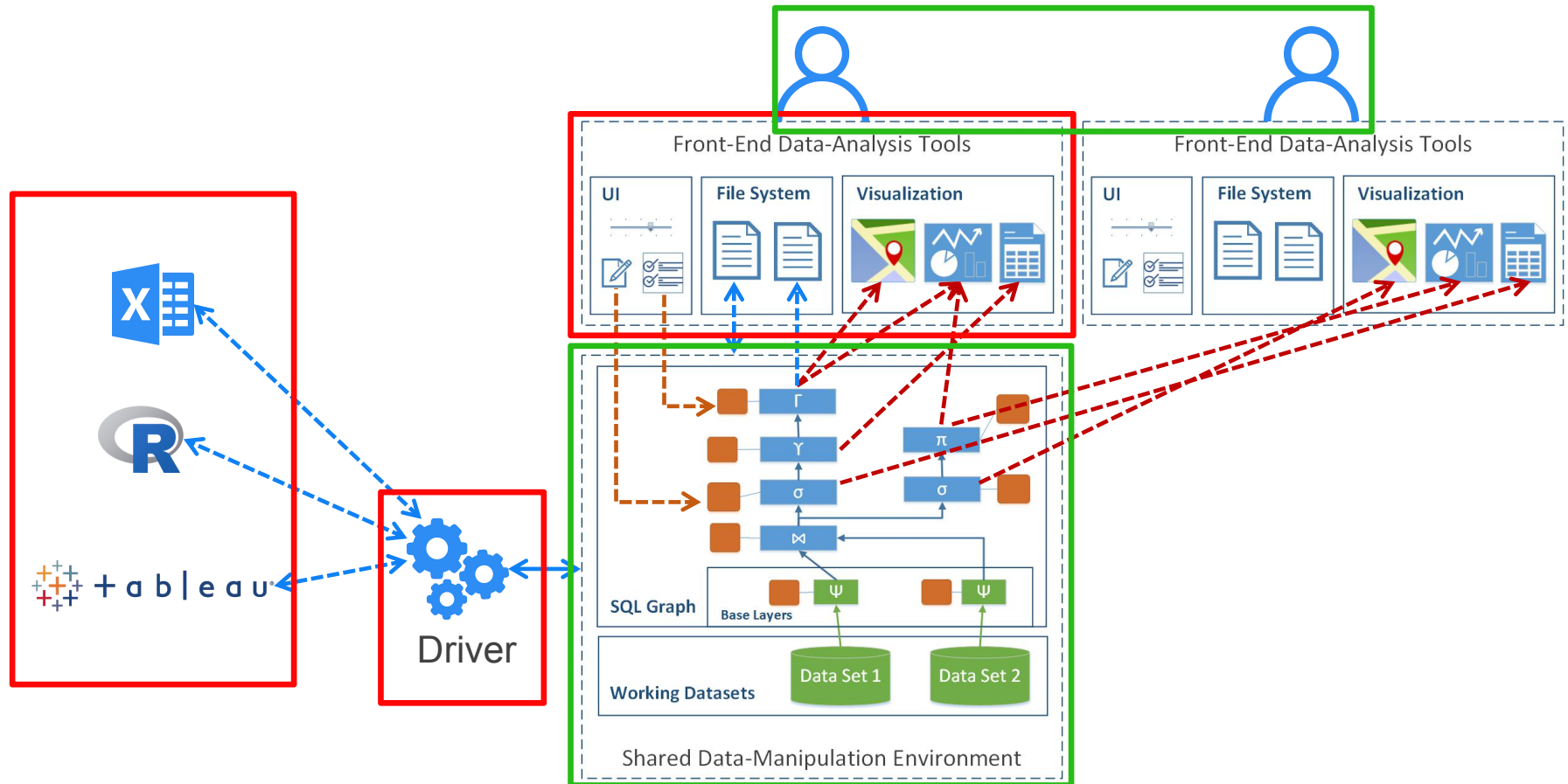  extremely tedious.

# Our Vision

# What Can We Do With jSQL$_e$?

# What Can We Do With jSQL$_e$?

# What Can We Do With jSQL$_e$?

# How Far Are We?

# Thank You