# llama Big Data Integration and Analysis

## Authors:
### Gaurav Kumar, Héctor Ugarte, Miguel Mármol, Tina Boroukhian
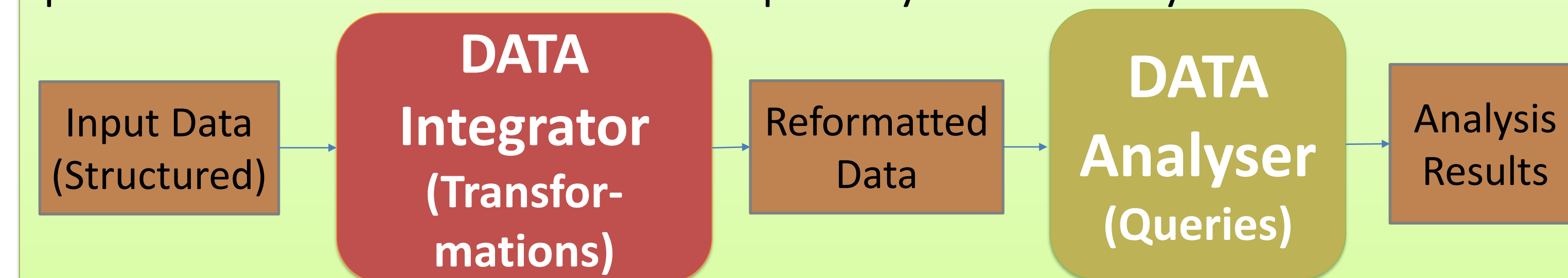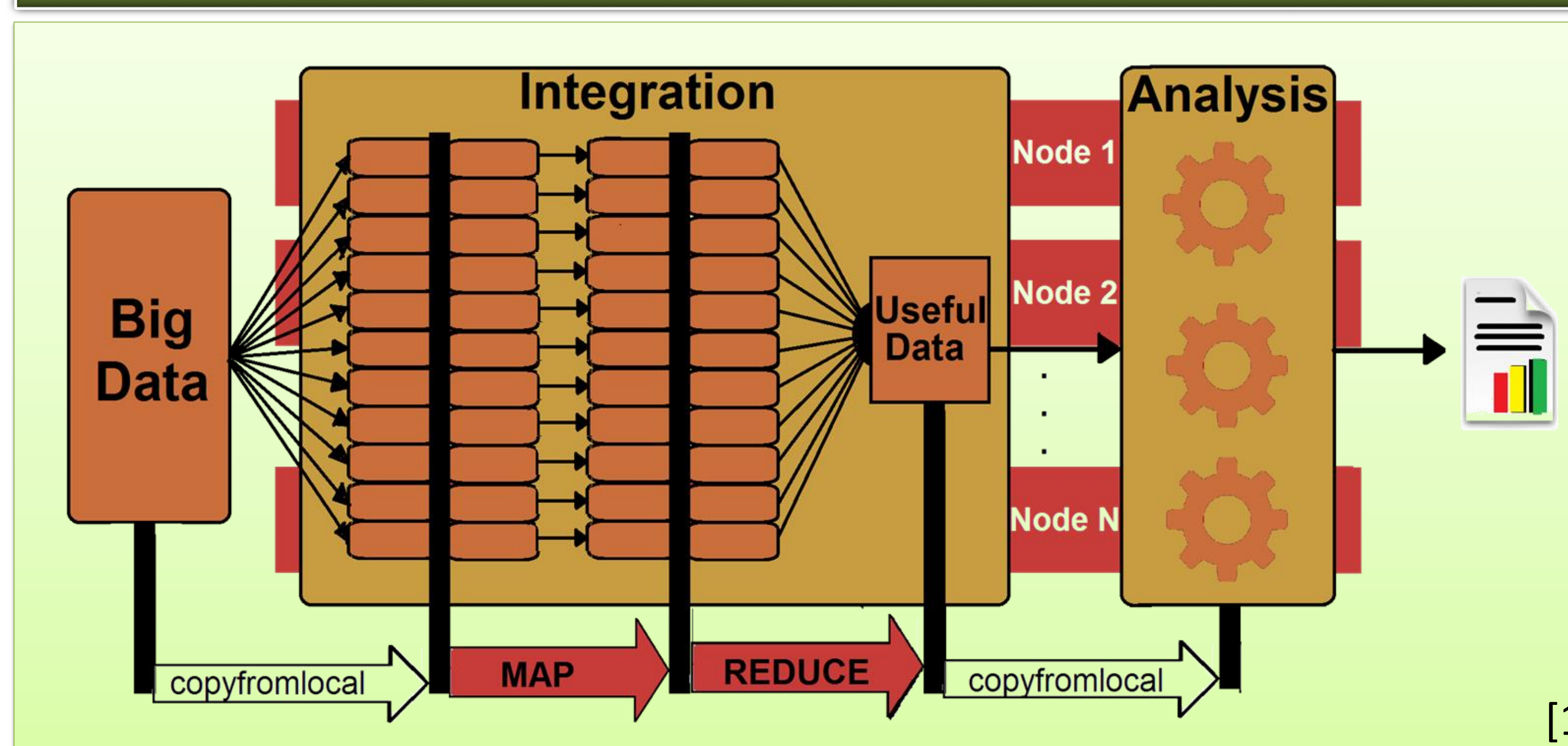### Summer Semester 2015

**BIG DATA**

## Project Overview

**Llama** is Data Integration and Analysis Platform made of two components:

**1. Data Integrator:** Llama first loads structured plain files. Then data loaded is cleaned, reformatted and filtered using a series of user-selected transformations. Integration jobs can be specified in two ways: either via (1) a graphical User Interface, or (2) a script written in **IJSL**, for **I**ntegration **J**ob **S**pecification **L**anguage. If the first method is used, at the end, the job can be exported as an IJSL script that can be used in a later run.

**2. Data Analyzer:** Once ready, the new data is analyzed by means of SQL queries. The results can be stored for possibly further analysis.
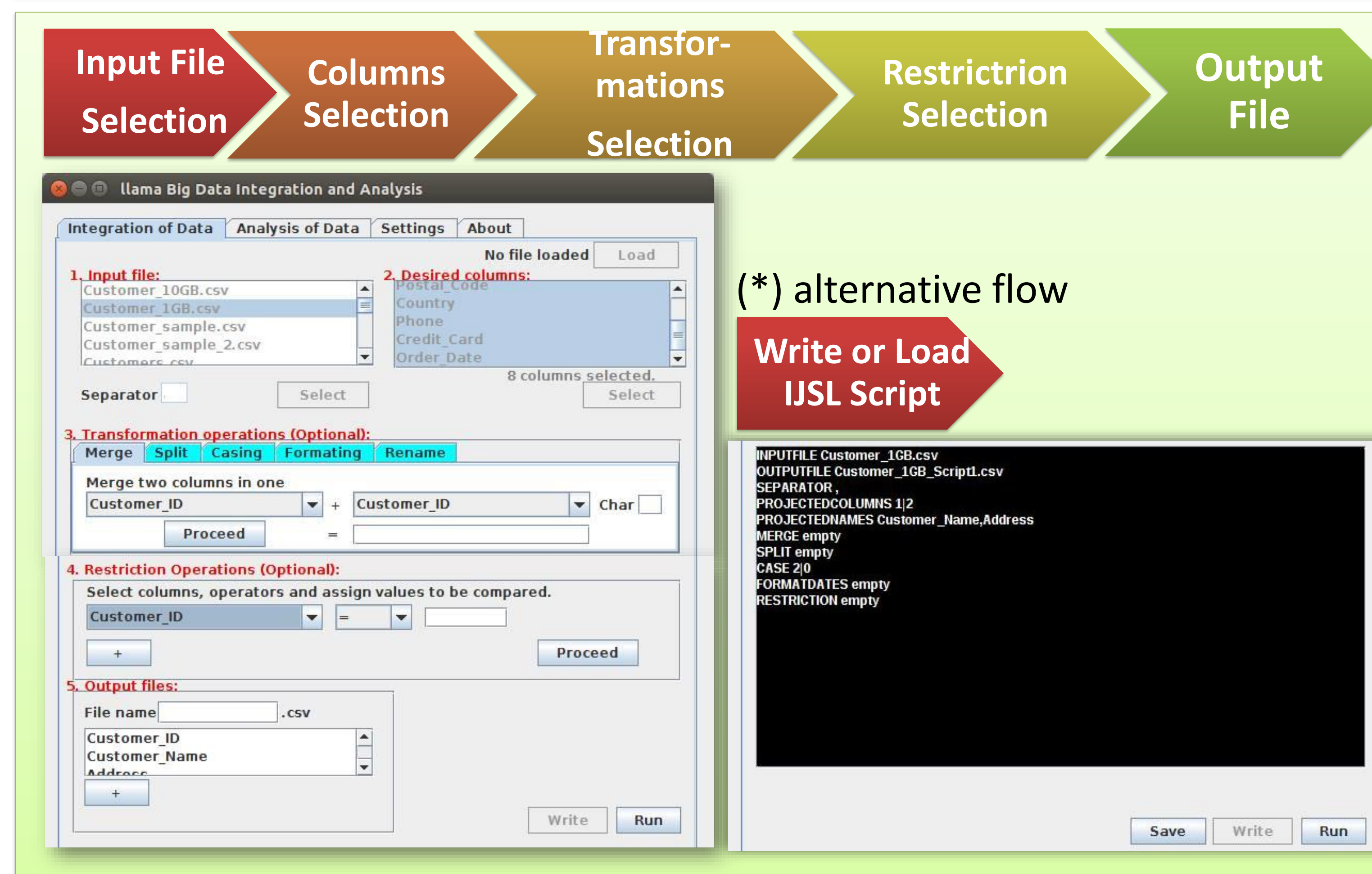


## System Architecture



[1]

## Implementation

**Apache Hadoop** is an open-source framework for distributed storage and processing of very large data sets on a of commodity hardware. [2]
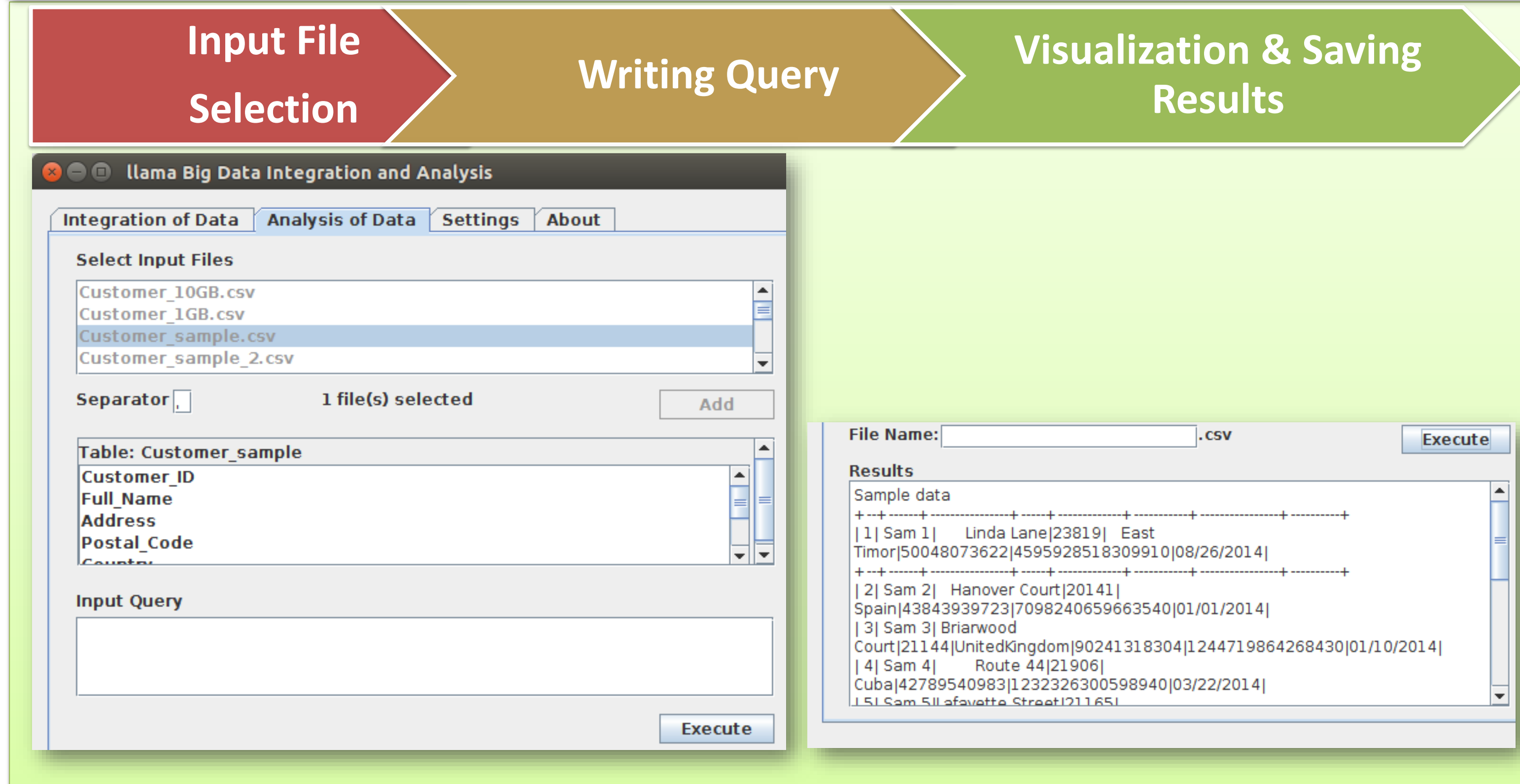
**Apache Spark** is an open-source cluster computing framework. It employs the concept of RDDs which are distributed units of data that resides primarily in memory, hence its high speed. [3]
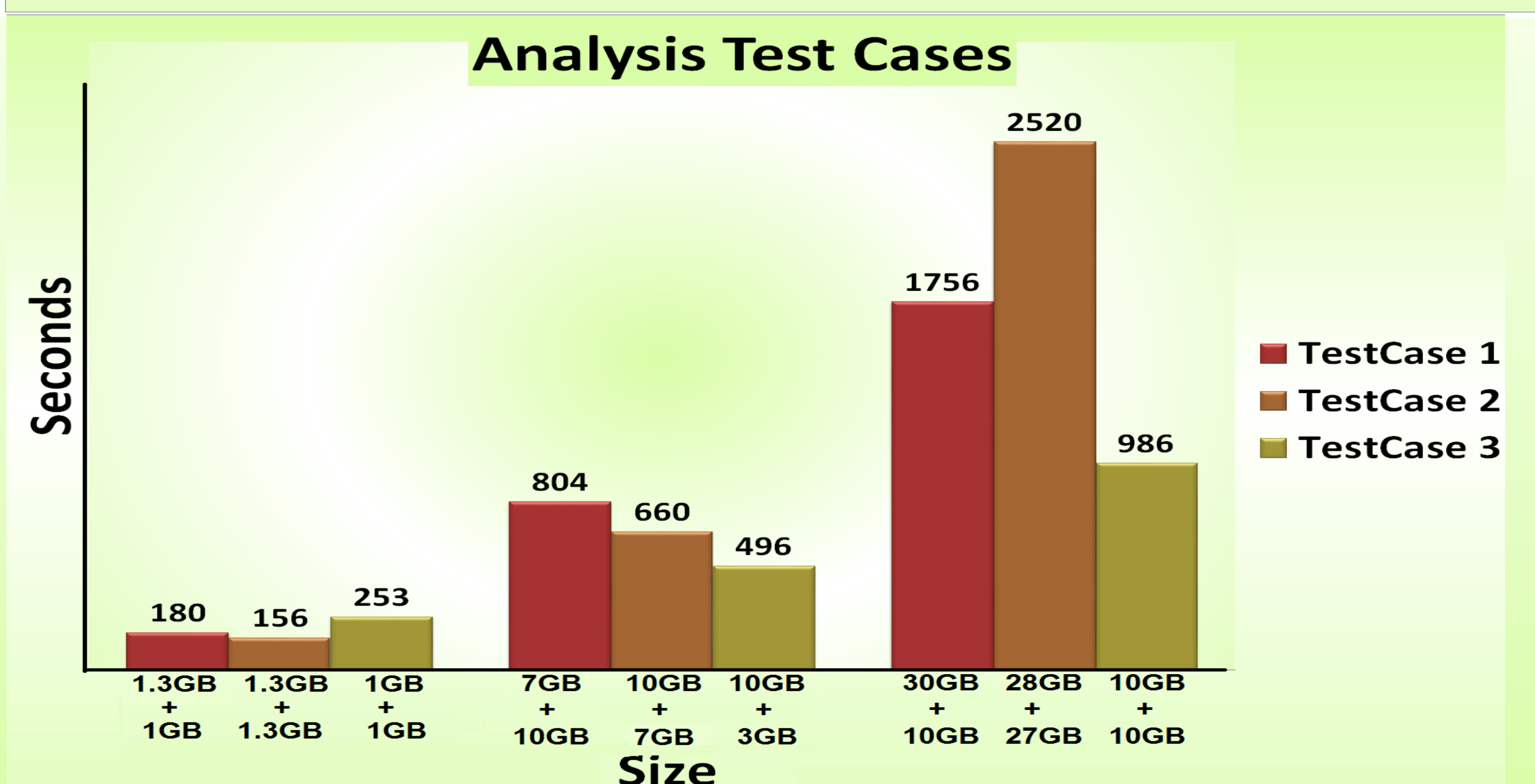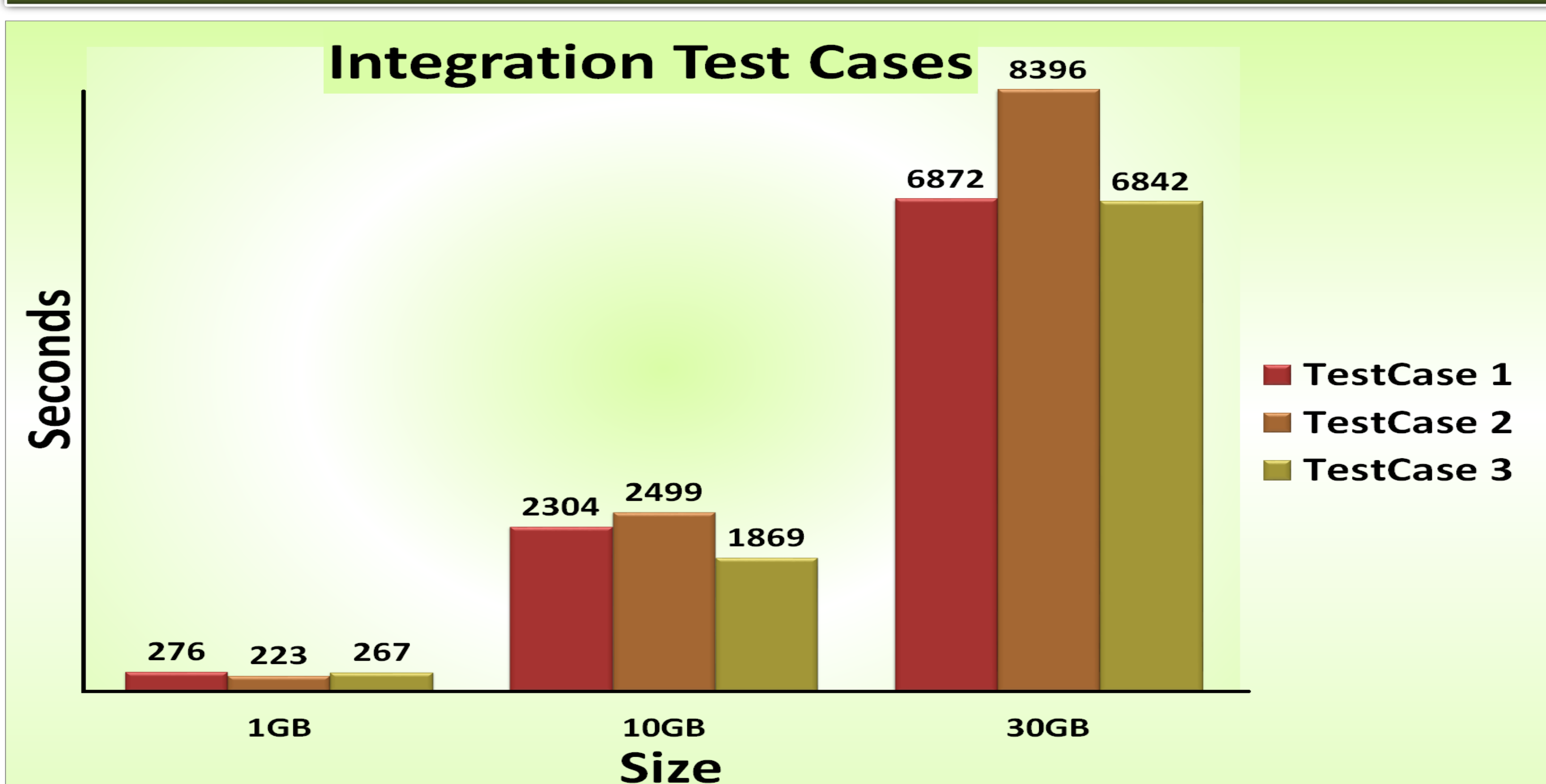
## Data Integrator



Input File Selection → Columns Selection → Transformations Selection → Restrictrion Selection → Output File

(*) alternative flow

Write or Load IJSL Script

## Data Analyser



Input File Selection → Writing Query → Visualization & Saving Results

## Evaluations

### Integration Test Cases



| | 1GB | 10GB | 30GB |
|---|---|---|---|
| TestCase 1 | 276 | 2304 | 6872 |
| TestCase 2 | 223 | 2499 | 8396 |
| TestCase 3 | 267 | 1869 | 6842 |

### Analysis Test Cases



| | 1.3GB+1GB | 1.3GB+1.3GB | 1GB+1GB | 7GB+10GB | 10GB+7GB | 10GB+3GB | 30GB+10GB | 28GB+27GB | 10GB+10GB |
|---|---|---|---|---|---|---|---|---|---|
| TestCase 1 | 180 | | | 804 | | | 1756 | | |
| TestCase 2 | | 156 | | | 660 | | | 2520 | |
| TestCase 3 | | | 253 | | | 496 | | | 986 |

## Live Demo

### A virtual machine hosted on EIS03

**Software installed:**
- Ubuntu 14.04 LTS 64bit
- 1GB RAM
- JDK 1.8
- Apache Hadoop 2.6
- Apache Spark 1.4

**Local Access:**
User: eis-user
Password: E1sbda2015

**Remote Access:**
Partner ID: 392368182
Password: E1sbda2015

References :
[1] http://www.glennklockwood.com/data-intensive/hadoop/mapreduce-workflow.png
[2] https://en.wikipedia.org/wiki/Apache_Hadoop
[3] https://en.wikipedia.org/wiki/Apache_Spark

**EIS Enterprise Information Systems**

universität **bonn**