# llama Big Data Integration and Analysis

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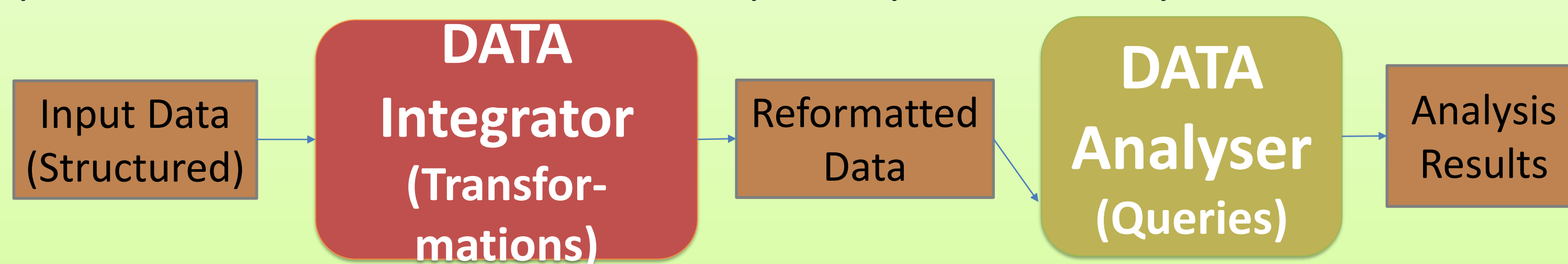### Summer Semester 2015

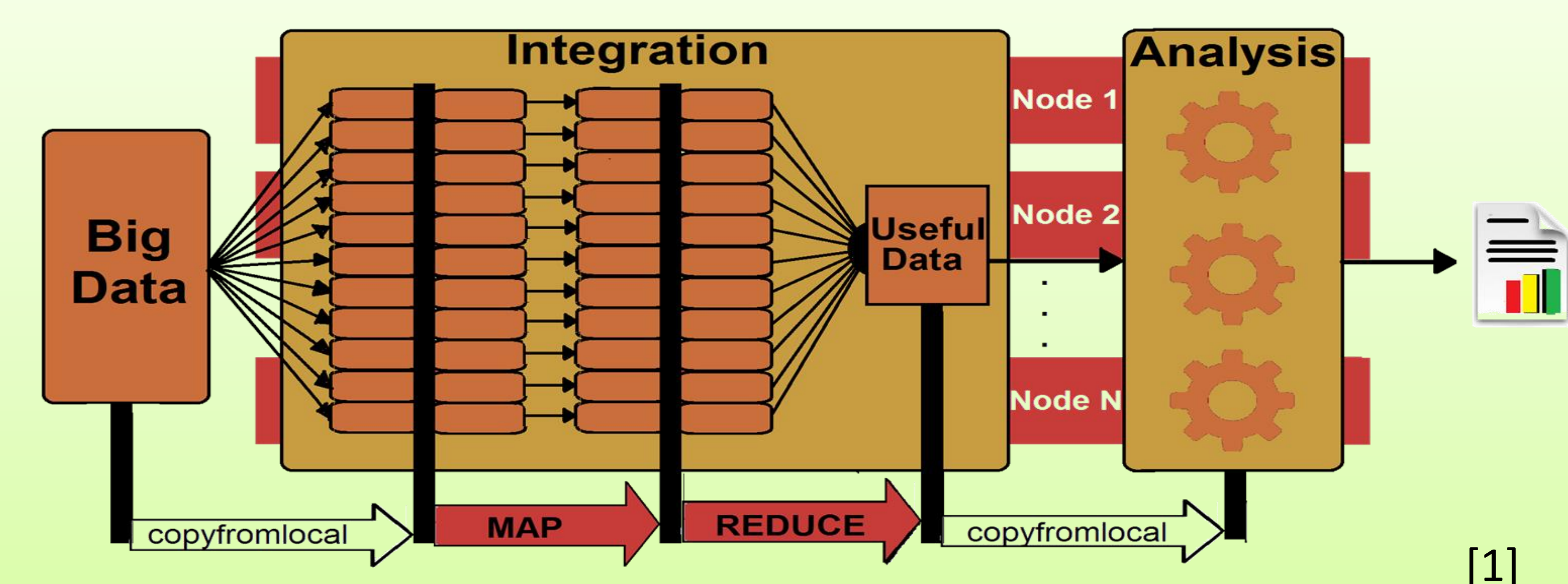**BIG DATA**

## 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.



Input Data (Structured) → DATA Integrator (Transformations) → Reformatted Data → DATA Analyser (Queries) → Analysis Results

## System Architecture



Big Data → copyfromlocal → MAP → REDUCE → copyfromlocal → Useful Data → Integration / Analysis → Node 1, Node 2, ... Node N

[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]
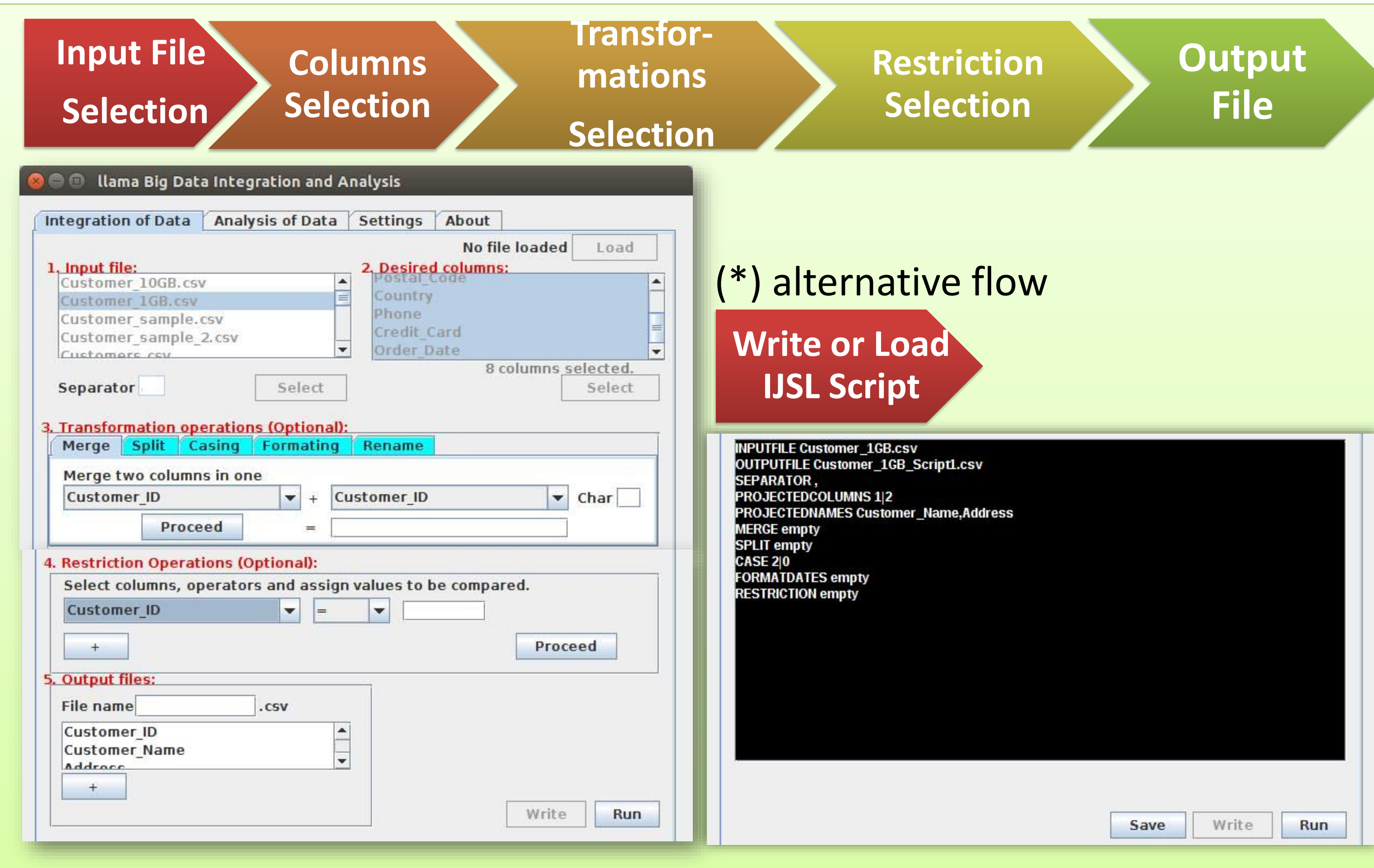
Tutorial for installation: http://pingax.com/install-apache-hadoop-ubuntu-cluster-setup/

**Apache Spark** is an open-source cluster computing framework. It employs the concept of RDD which are distributed units of data that resides primarily in memory, hence its high speed. [3]
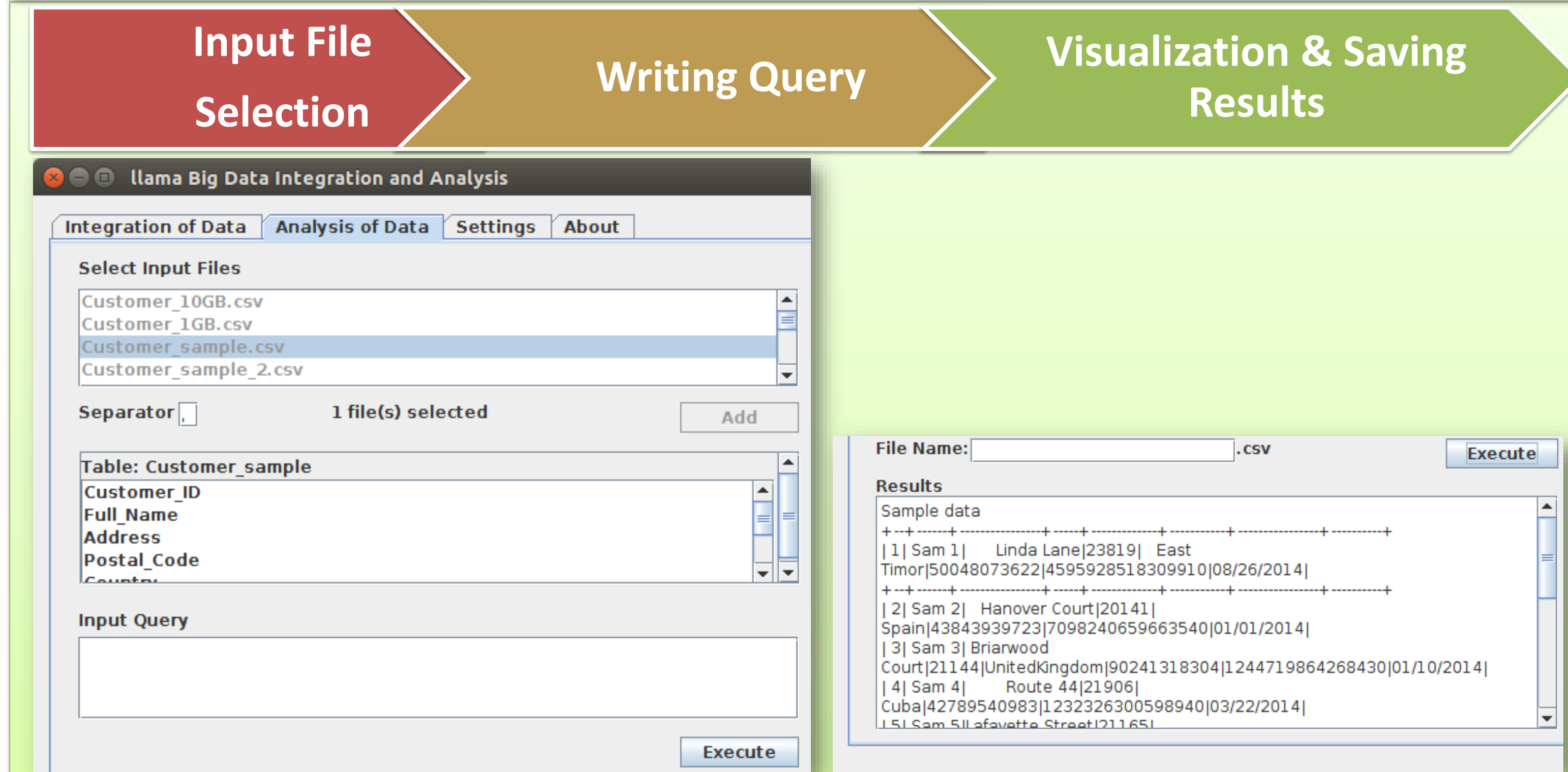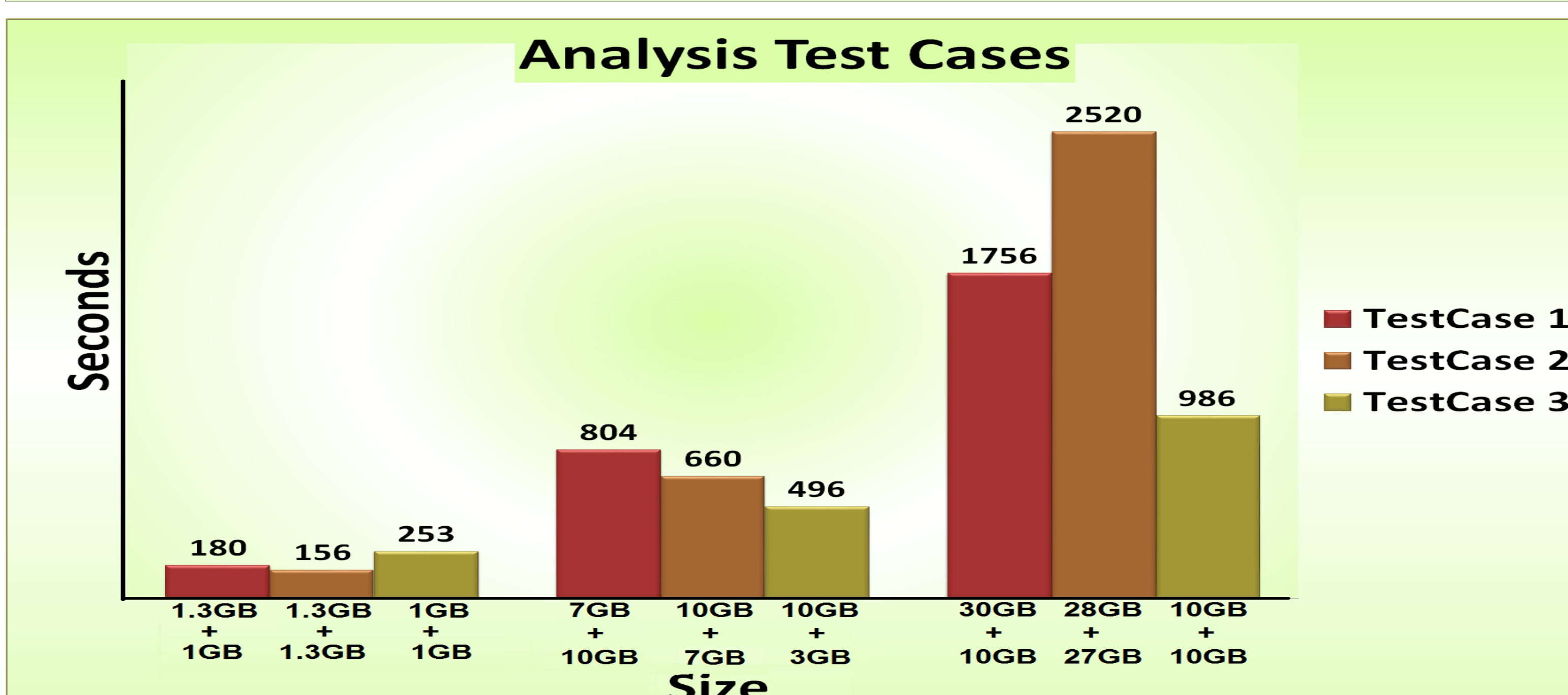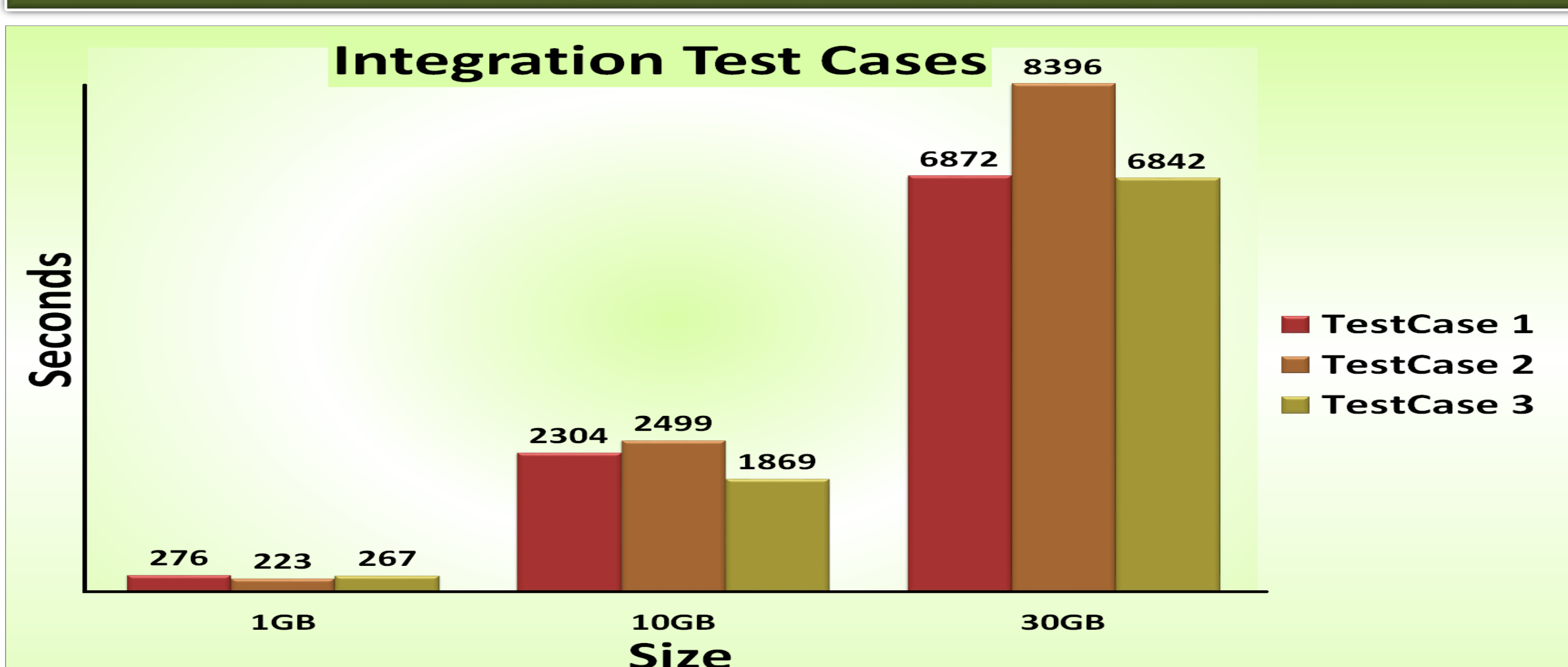
Tutorial for installation: http://www.trongkhoanguyen.com/2014/11/how-to-install-apache-spark-121-in.html

## Data Integrator

Input File Selection → Columns Selection → Transformations Selection → Restriction Selection → Output File

(*) alternative flow

Write or Load IJSL Script



## Data Analyser

Input File Selection → Writing Query → Visualization & Saving Results



## Evaluations



### Integration Test Cases

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



### Analysis Test Cases

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

## Live Demo

### Two virtual machines hosted on EIS02

**Local Access Master:**
User: eis-groupb-master
Password: EIS2015

**Hadoop User:**
User: hduser
Password: hduser

**Remote Access:**
Partner ID: 572120720
Password: 268vxh

**Local Access Node:**
User: eis-groupb-slave
Password: EIS2015

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

**Remote Access Node:**
Partner ID: 572193139
Password: 4d2g1h

Github Link:
[1] **User Manual:** https://docs.google.com/document/d/1Rhdg8mlKdiJpZoFBIyfrao1-czLz4FkjaBdnLxr8V60/edit
[2] **Technical Report:** https://docs.google.com/document/d/1kNYN4YxueQ0UbKZSeWWSwx5JEsuJFOX08XeEkN7a8S0/edit

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**