



Llama Big Data Integration and Analysis

Authors:

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

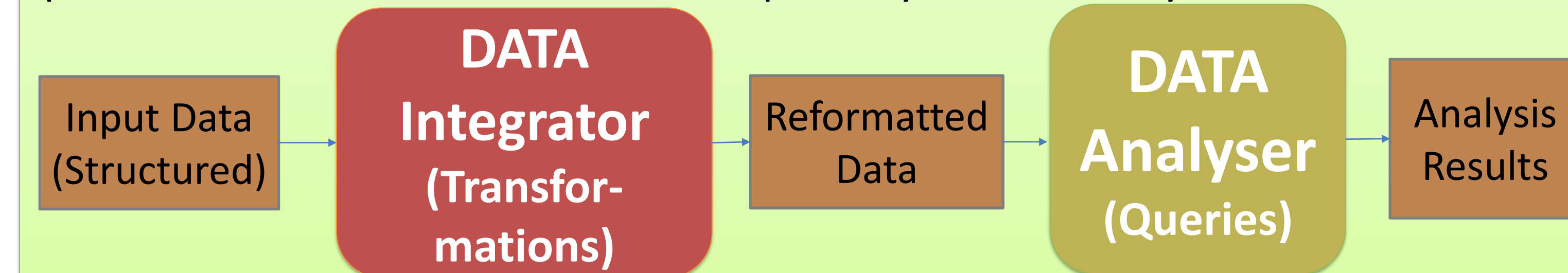
BIG
DATA

Project Overview

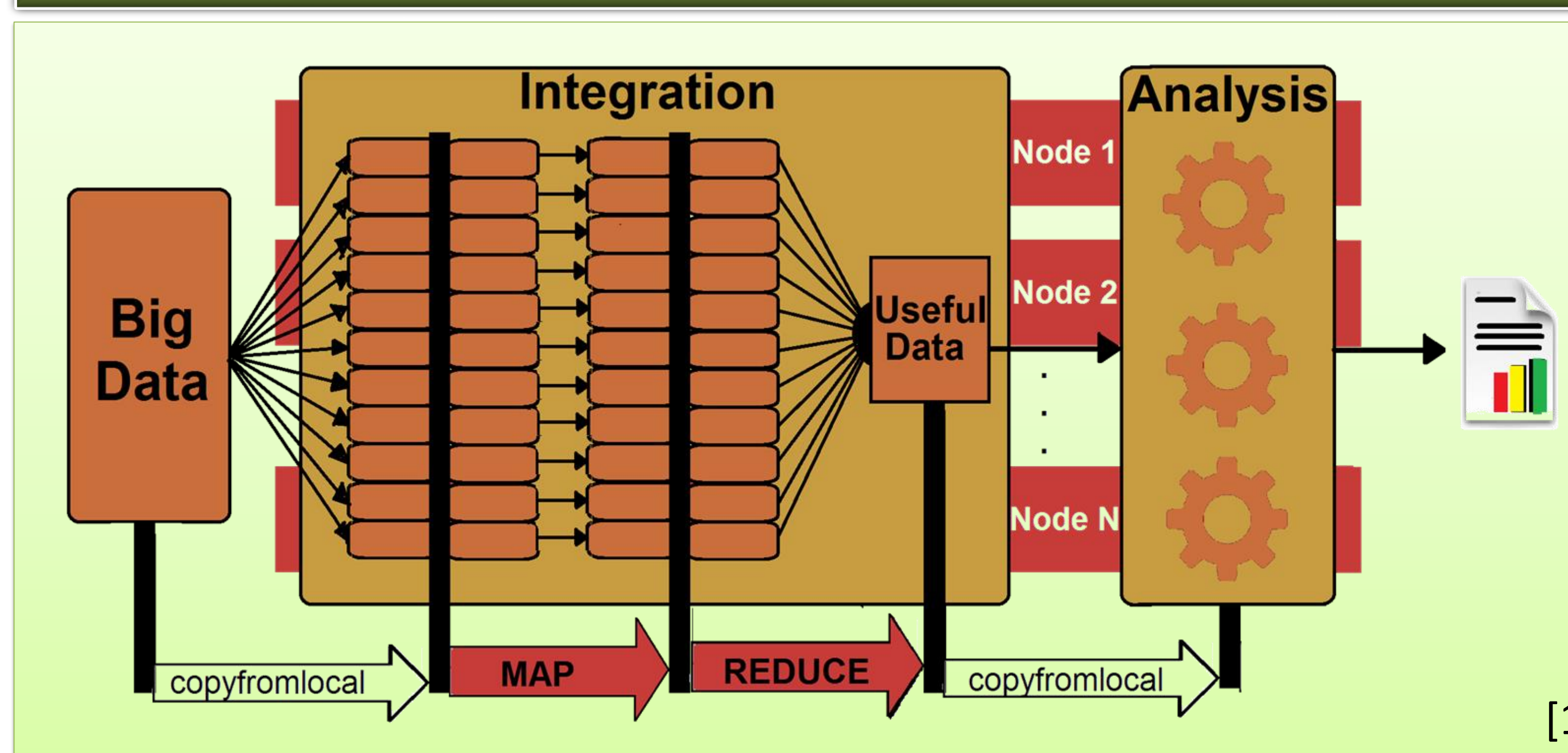
Llama is Data Integration and Analysis Platform made of two components:

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 Integration Job Specification Language. 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.

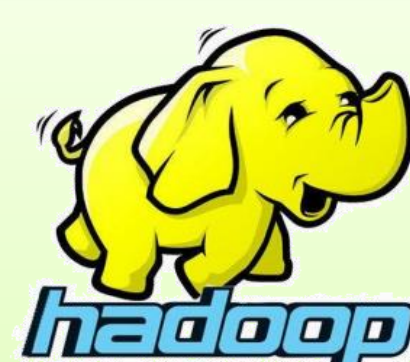


System Architecture



Implementation

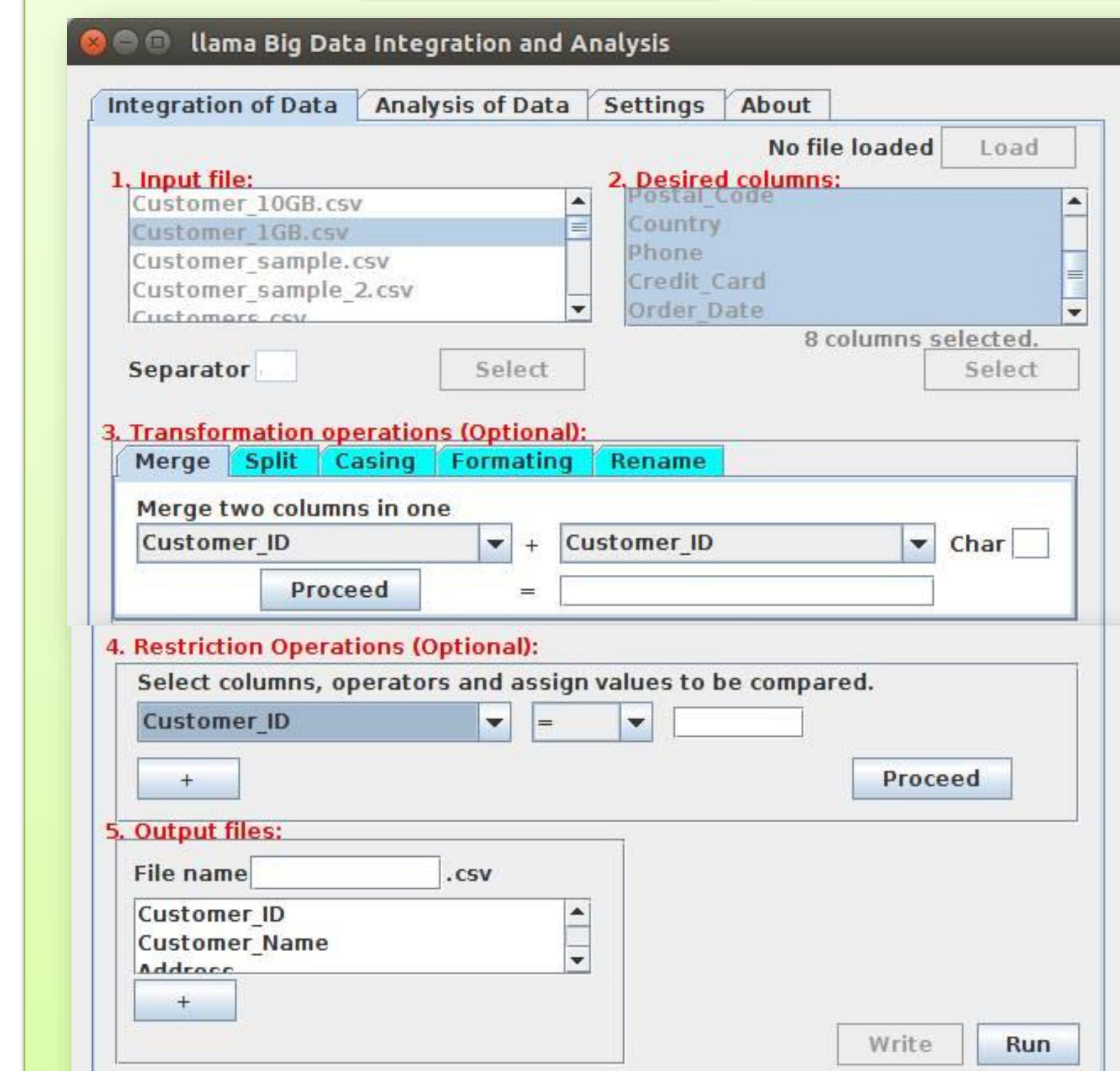
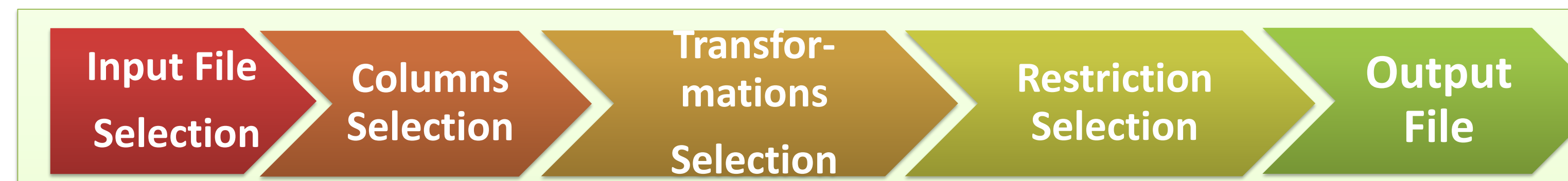
Apache Hadoop is an open-source framework for distributed storage and processing of very large data sets on a of commodity hardware. [2]



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]

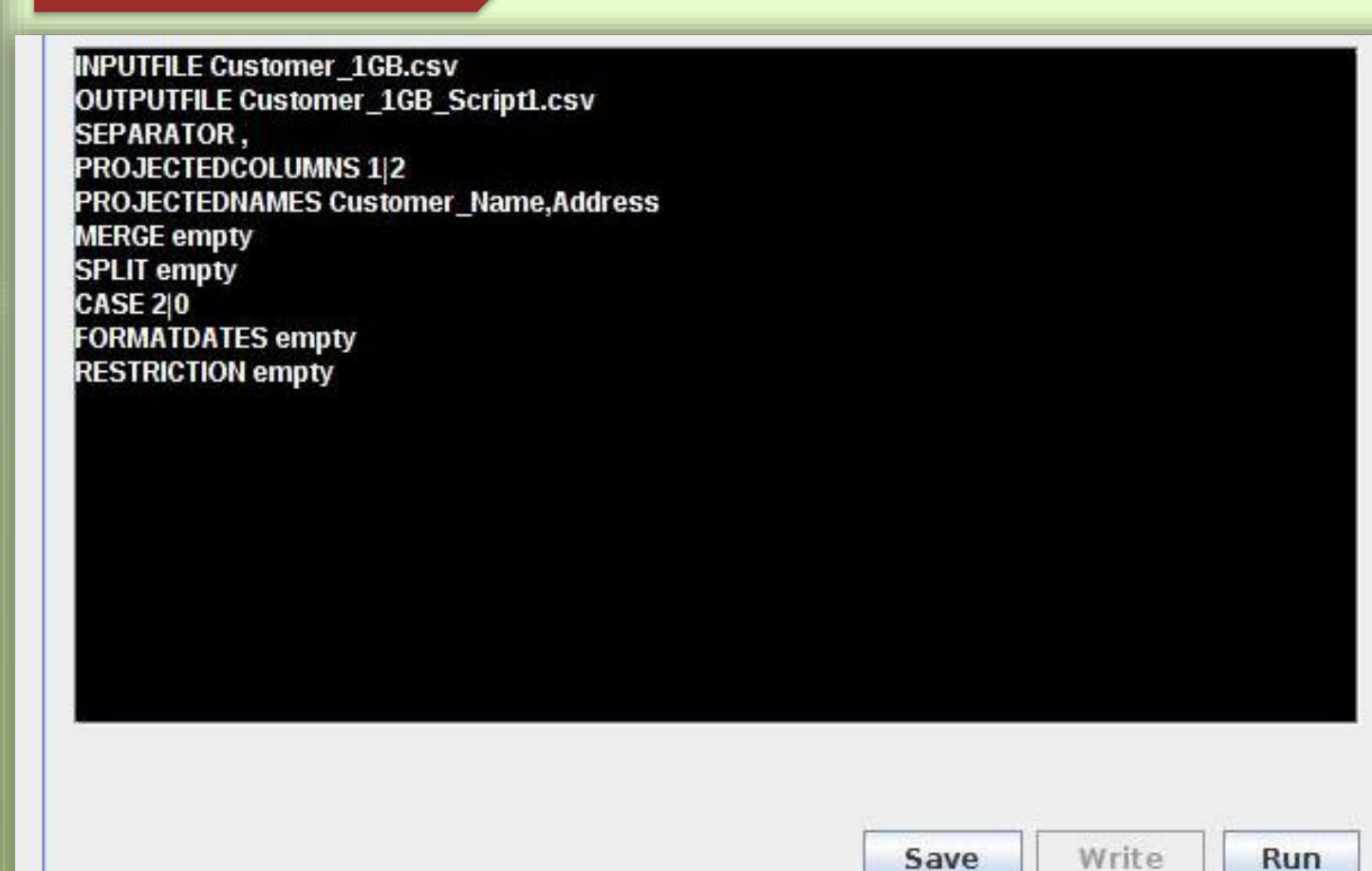


Data Integrator

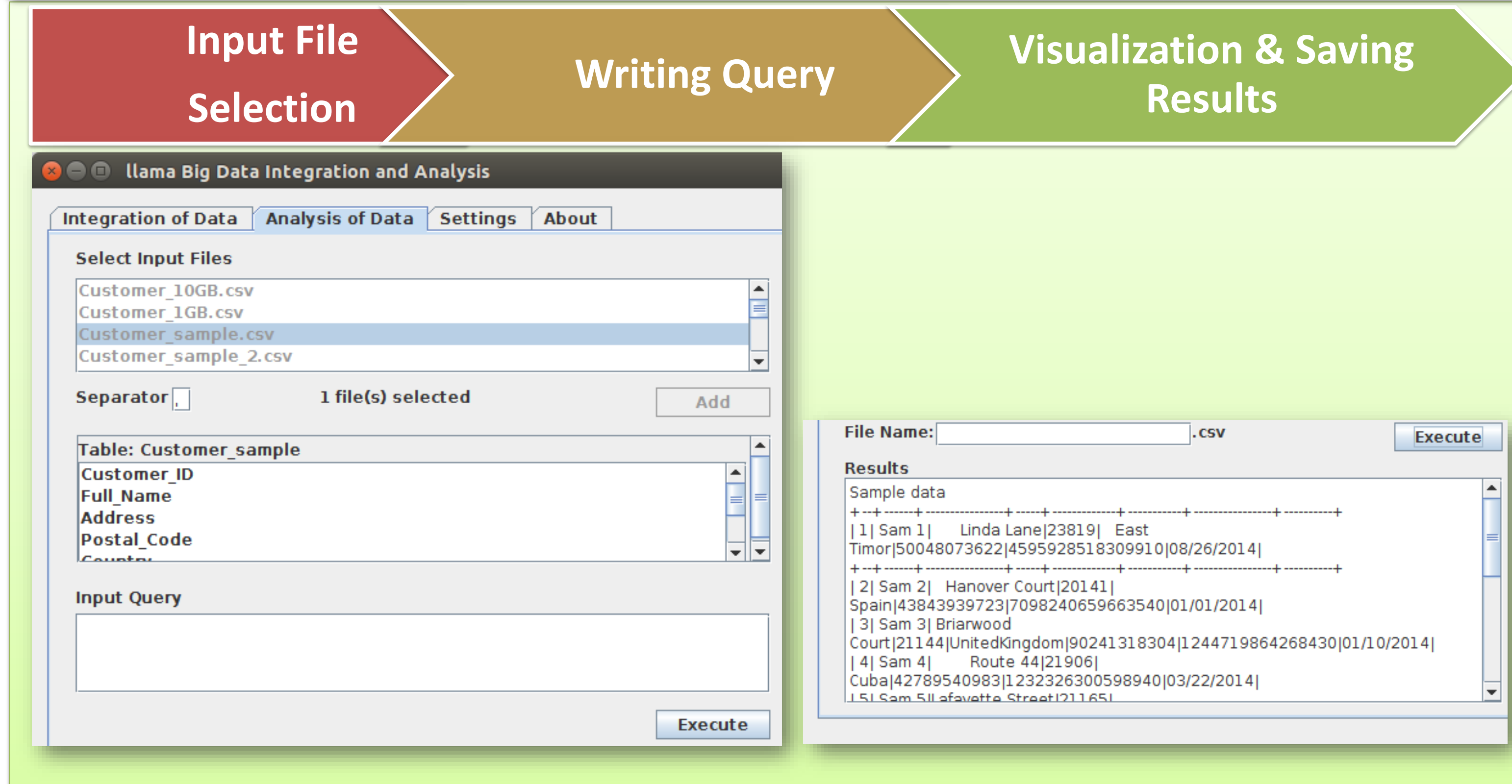


(*) alternative flow

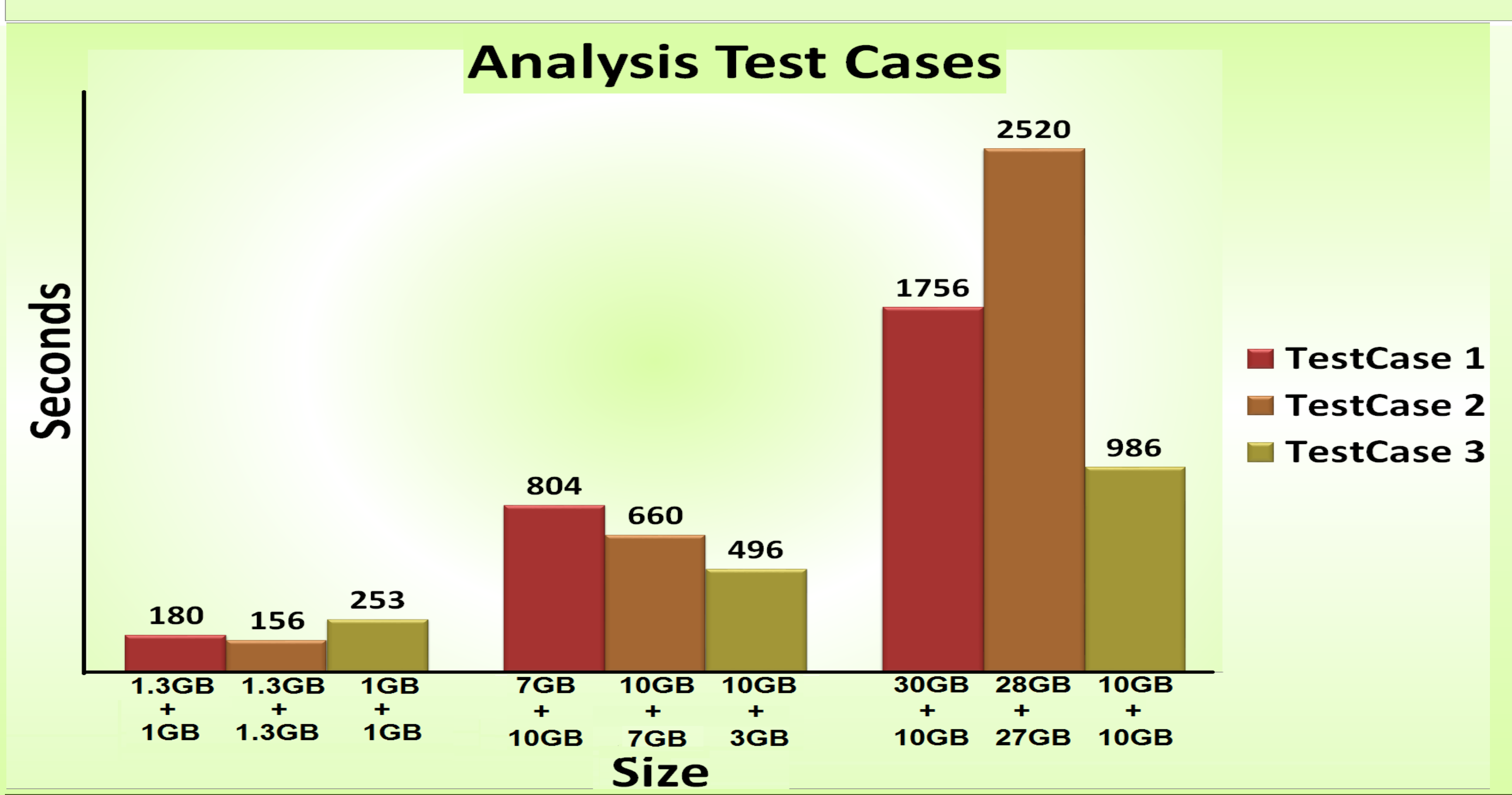
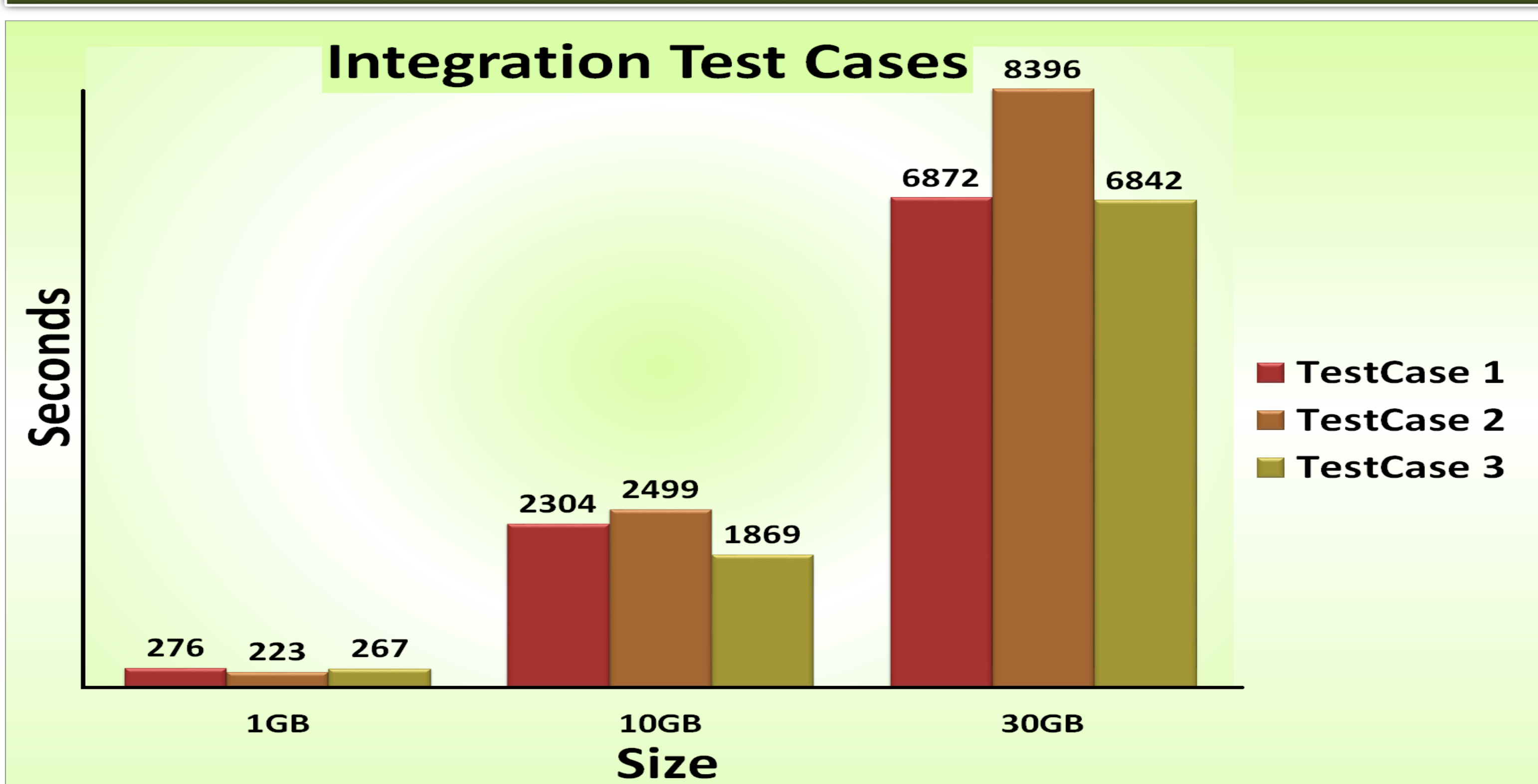
Write or Load
IJSL Script



Data Analyser



Evaluations



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 Master: Partner ID: 57210720 Password: 268vxh
Local Access Node: User: eis-groupb-slave Password: EIS2015	Software installed: <ul style="list-style-type: none">Ubuntu 14.04 LTS 64bitJDK 1.8Apache Hadoop 2.6Apache Spark 1.4	Remote Access Node: Partner ID: 572193139 Password: 4d2g1h

References :

- <http://www.glennklockwood.com/data-intensive/hadoop/mapreduce-workflow.png>
- https://en.wikipedia.org/wiki/Apache_Hadoop
- https://en.wikipedia.org/wiki/Apache_Spark