# Data Munging and Analysis for Scientific Applications

Raminder Singh
Science Gateways Group
Indiana University, Bloomington
raminder@apache.org

# Overview

- Evaluate the Apache Big data tools
- Understand the execution patterns of Analysis applications
- Solutions using Airavata
- Build a gateways solution with HPC and Big Data requirements

# Hadoop 2 Ecosystem

**Applications Run Natively IN Hadoop**

| BATCH (MapReduce) | INTERACTIVE (Tez) | ONLINE (HBase) | STREAMING (Storm, S4,...) | GRAPH (Giraph) | IN-MEMORY (Spark) | HPC MPI (OpenMPI) | OTHER (Search) (Weave...) |

**YARN** (Cluster Resource Management)

**HDFS2** (Redundant, Reliable Storage)

# **Motivation to explore**

- Heterogeneous data
- Data Munging (parsing, scraping, formatting data)
- Visualization or Analyze
- Preservation of data

# Analysis Applications

- Behavior Tracking - medical

- Situational Awareness - weather

- Time Series Data -Patient monitoring, weather data to help farmers

- Resource consumption Monitoring - Smart grid

- Process optimization

# Scientific applications Data Types

Observational Data – uncontrolled events happen and we record data about them.

Examples include astronomy, earth observation, geophysics, medicine, commerce, social data, the internet of things.

Experimental Data – we design controlled events for the purpose of recording data about them.

Examples include particle physics, photon sources, neutron sources, bioinformatics, product development.

Simulation Data – we create a model, simulate something, and record the resulting data.

Examples include weather & climate, nuclear & fusion energy, high-energy physics, materials, chemistry, biology, fluid dynamics.

# What is Science Gateway?

- Community portal or desktop tools
- Common science theme
- Collaborative environment

The Ultrascan science gateway supports high performance computing analysis of biophysics experiments using XSEDE, Juelich, and campus clusters.

We help build gateways for labs or facilities.

Launch analysis and monitor through a browser

Desktop analysis tools

# BioVLAB

# Airavata

# Value of using Airavata

- Enable collection of resources
- Application centric not compute centric
- Meta workflow to enable set of applications

# Use-case for Data Analysis

- TextRWeb: Large Scale Text Analytics with R on the web

Collaborator: Hui Zhang, Data Scientist at Indiana University

**TextRWeb web service**

Document-centric API

Restricted analytics on TDM

Result

Returned results

```
library(multicore)
library(tm)

# the directory where user data resides,
# passed from TextRWeb front end
documentHome <- "../../R/library/tm/texts/htrc_volumes"
userCorpus <- Corpus(DirSource(documentHome, encoding = "UTF-8"),
    readerControl = list(language = "en")))

################################################
# user defined document-centric function
################################################

wrapper <- function(document) {
    # define custom tokenizer
    strsplit_space_tokenizer <- function(x) unlist(strsplit(x,
        "[[:space:]]+"))

    # minimal frequency threshold
    minFreqThreshold <- 2

    # control list
    ctrl <- list(tolower = TRUE, tokenize = strsplit_space_tokenizer,
        removePunctuation = list(preserve_intra_word_dashes = TRUE),
        removeNumbers = TRUE,
        stopwords = stopwords("english"),
        stemming = TRUE,
        wordLengths = c(minFreqThreshold, Inf))

    # generate term frequency
    termFrequency <- termFreq(document, control = ctrl)

    # return termFrequency
    termFrequency

}

################################################
# apply wrapper function to each
# in parallel by 'mclapply' document
################################################
numCores <- 8
result <- mclapply(userCorpus, wrapper, mc.preschedule = TRUE,
    mc.cores = numCores)
```
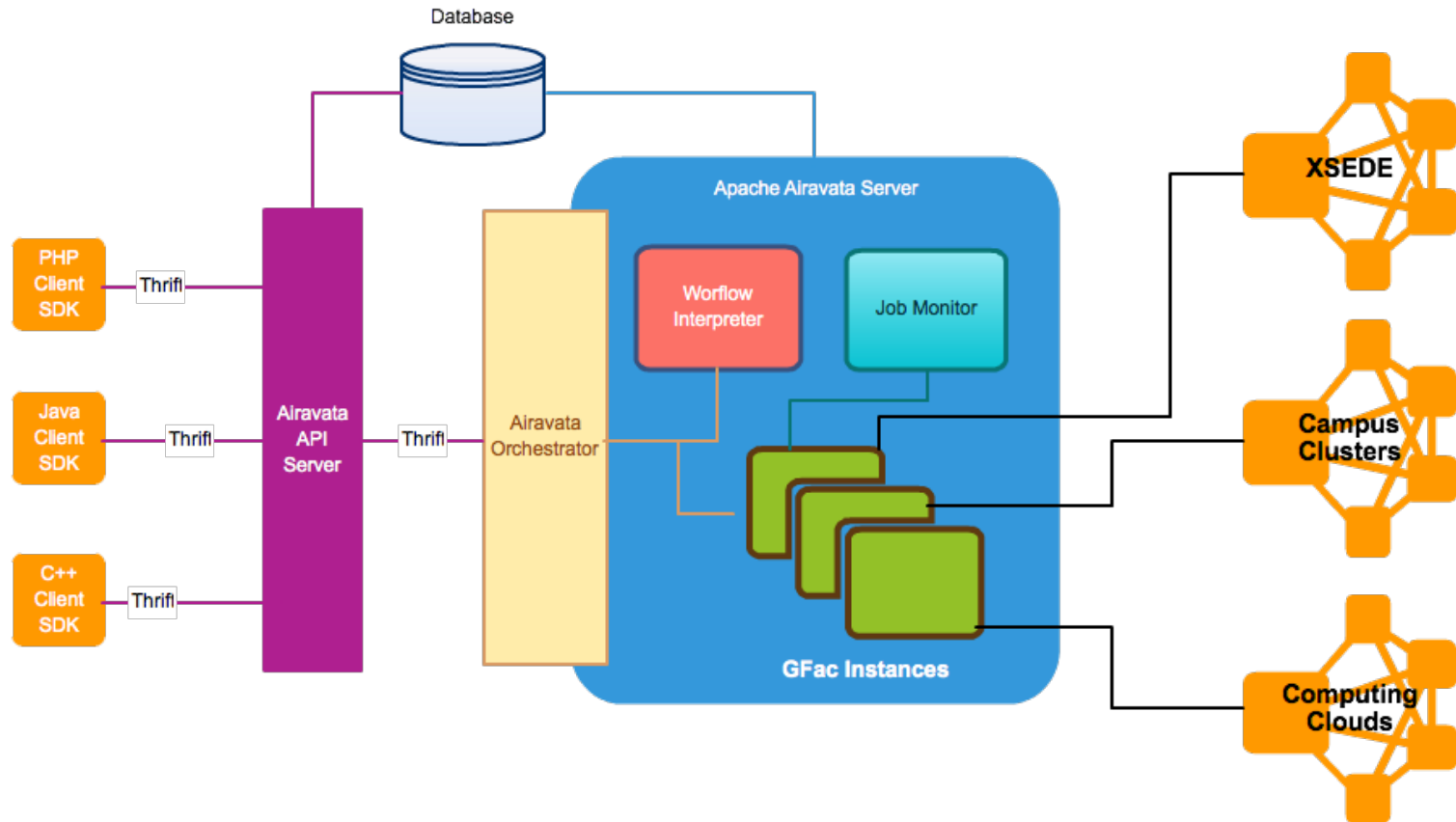
**Complete code generation**

**Load balancers**

**HTRC data services**

**DNS**

**Worker hosts**

**Cloud/HPC**

1  2  3  4  5  6  7  8  9

**Goals for R on the web project**

- Run large scale text analysis using parallel R.
- Hide computational complexity with user interfaces
- Support interactive text analysis
- Support iterative text mining

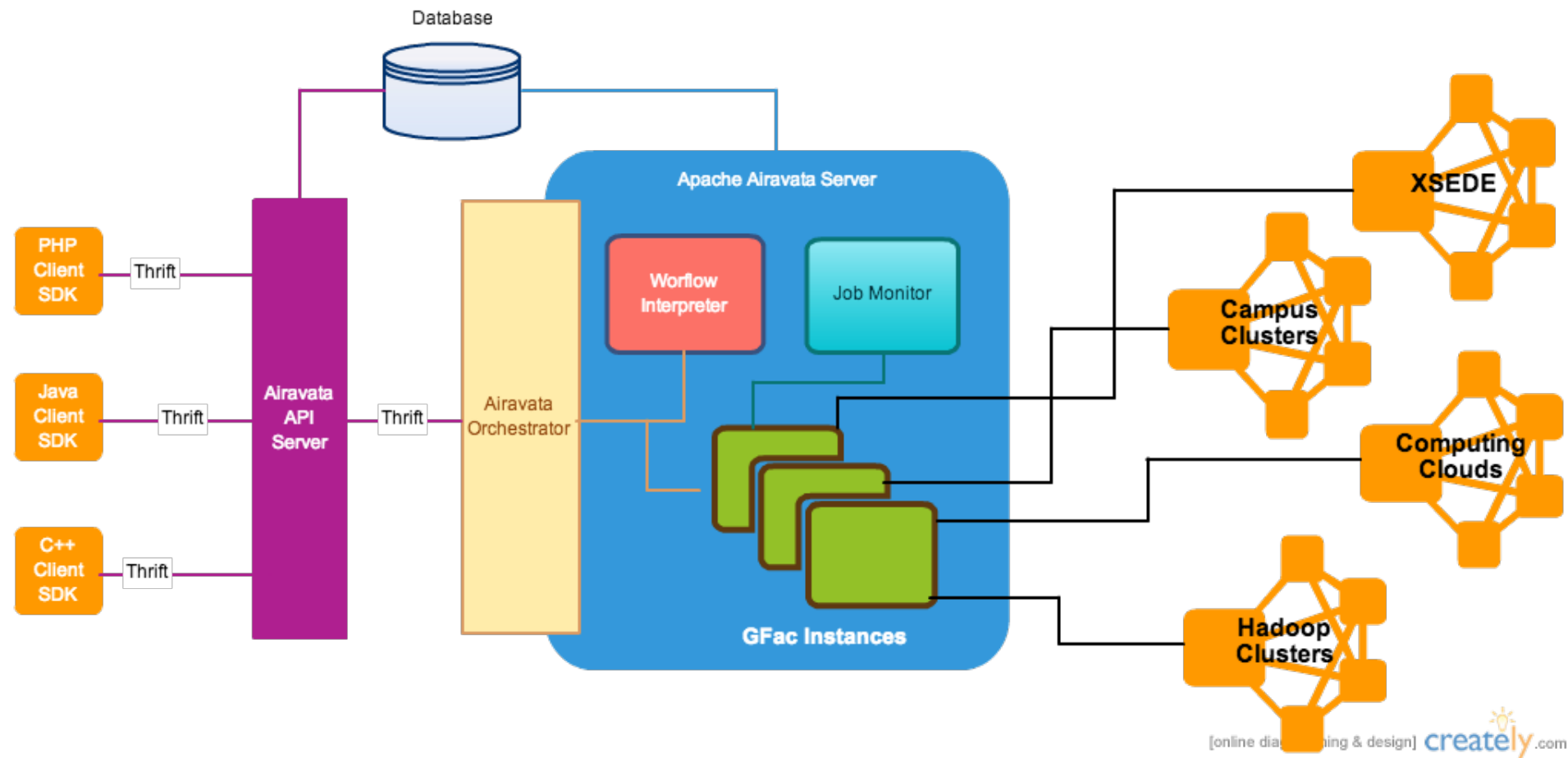# TextR Solution Diagram

# Current Hadoop Integration

# Future Work

- Integrate TextRWeb with Apache Spark
- Explore SparkR [1]
- Develop Apache Thrift interfaces for TextRWeb server
- Integrate with Apache Airavata for HPC job.
- Explore workflow DAGs for Text Analysis
- Keep updated with product offering like Stratosphere

1. https://github.com/amplab-extras/SparkR-pkg

# Apache Spark

- In Memory computations

- Machine learning library (MLLib)

- graph engine (GraphX)

- Streaming analytics engine (Spark Streaming)

- Fast interactive query tool (Shark).

- Use Lineage data for fault tolerance

  - Tracking the data path

Database

Apache Airavata Server

PHP
Client
SDK

Thrift

Java
Client
SDK

Thrift

C++
Client
SDK

Thrift

Airavata
API
Server

Thrift

Airavata
Orchestrator

Worflow
Interpreter

Job Monitor

GFac Instances

XSEDE

Campus
Clusters

Computing
Clouds

Hadoop
Clusters

# Conclusion

- Value added for the scientific communities
- Value for Apache Big Data Suite

# Q & A

## airavata.apache.org

**Subscribe: users-subscribe@airavata.apache.org**
**Subscribe: dev-subscribe@airavata.apache.org**
**Subscribe: architecture-subscribe@airavata.apache.org**

# Thanks You!