# Map Reduce for Algorithmic Trading

Akshaan Kakar, Alex Berard
*Dept. of Computer Science*
*Columbia University*
*Email: ak3808@columbia.edu, alice.berard@columbia.edu*

*Abstract*—**Algorithmic Trading is an extremely competitive sector of financial markets. Developing trading algorithms involves the pivotal step of backtesting, where the performance of the algorithm is validated against large amounts of historical securities pricing data. These time series require large amounts of computational resources for storage, processing and visualization. In this paper, we describe the implementation of an algorithmic trading backtest engine that uses "Big Data" tools as a backbone to backtests for trading algorithms in an efficient and scalable fashion. We also show how this system can be extended to support live trading as well as a plethora of other features.**

*Keywords*-**Algorithms, Finance, Trading, Big Data, Hadoop, Spark, Backtesting**

## I. INTRODUCTION

With the computerization of order flows in financial markets, traders were afforded the ability to have computers place buy and sell orders on securities according to pre-defined strategies. The advent of such Algorithmic Trading strategies eventually developed into an extremely competitive sector of financial markets. The development of trading algorithms is now a formalized process that involves intricate research and mathematical models. Although a significant portion of this development effort is invested in the underlying mathematical theory, backtesting of algorithms is the most pivotal step in the process. Backtesting is the phase in which a trading algorithm is validated against large amounts of historical pricing data. Since trading algorithms may place orders extremely frequently and involve multiple securities in a single portfolio, large amounts of data must be stored, read and processed in order to run these tests. Moreover, the result must be visualized so that developers can gauge the performance of their algorithms quickly and conveniently.

There is an abundance of well developed and well maintained tools which have been built expressly to grapple with the large amounts of data that have now become commonplace. These tools leverage novel algorithmic and system level techniques to to deal with the space and time bottlenecks that come with large datasets. Since the problem of backtesting algorithmic strategies is inherently a 'Big Data' shaped problem, we decided to use a combination of such tools in order to to build an efficient and flexible system for the task.

## II. RELATED WORKS

We modeled our implementation of the backtesting engine with some other tools and works in mind. One of the prominent platforms that we sought to emulate is Quantopian (CITE). The Quantopian platform allows in-browser implementation of trading algorithms in Python. It also allows the use of internal APIs as well as most external python modules for statistics, numerical algorithms etc. Although we have not extended our system to allow in-browser coding, we have also used python as the language of choice for trading algorithm implementation. Python is readable, flexible and popular, making it a good choice for algorithmic trading; a field where not everyone is well versed with computer programming. We also sought to emulate the charting and metric calculations performed by the Quantopian platform. The benchmark as well as the returns time series are plotted to provide a lucid understanding of the algorithm performance, without confusing the user with unnecessary detail. Also, important metrics such as the Sharpe Ratio, Maximum drawdown etc. are provided to the user for further insight.

## III. SYSTEM OVERVIEW

The main objectives in our implementation were to keep the platform efficient and easily extensible. In order to achieve these goals, we partitioned the system into independently functioning subsystems which we then integrated to realize the final backtesting engine. For the purpose of demonstration, we utilized the free historical pricing data from QuantCode (CITE). The dataset as well as the subsystems are described in detail in the following subsections.

### A. QuantQuote dataset

We decided to use the QuantQuote free historical stock price data for the purpose of demonstration. This data consists of daily stock tick data for the 500 symbols that are listed on the Standard and Poor's 500 Index from 1998 to present. Although we used the dataset throughout the design and testing process, we structured the backtesting engine such that it is dataset agnostic. More specifically, we designed our subsystems so that they do not impose any tight constraints on input data formatting. Further descriptions

of input data specifications are provided in the following subsection.

### B. Hadoop Data Warehouse

Apache Hadoop (CITE) is a distributed framework that is designed for the storage and processing of large datasets. Hadoop is especially well suited to read-only, batch accessing of large amounts of data. Our system only requires the reading of financial time series data, without frequent writes or editing, making Hadoop a perfect choice for the data warehouse subsystem. We stored the finance symbol pricing data in the Hadoop Distributed File System (HDFS), which can be configured to run on any number of machines without a single point of failure. We structured the file system so that the time series data for each symbol were stored in a separate file. This simple, flat structure makes accessing the data extremely simple, since for a particular symbol, only a single file must be accessed. Each file was named after the symbol, to facilitate programatic access using symbol names directly, without the need for any lookup or translation. In the case that the data repository needs to be updated with newer pricing data, only a single file need be written to. This simple organization does not compromise on efficiency since Hadoop is designed to work well with a relatively small number of very large files, as opposed to small fragments of a larger dataset. A directory listing showing the file structure is shown in Figure 1.

Within each file, the time series were stored in a comma separated value (CSV) format, with time stamps and price values as the fields in each row. The CSV format is simple and widely used, making our system largely data agnostic. Any data in CSV format with the appropriate fields can be used with our data warehouse system since our engine does not impose and other restrictions on data formatting.

### C. Spark Algorithm Processing Engine

Apache Spark is a tool that was developed for processing general large-scale data efficiently (CITE). Spark come with support for multiple languages and platforms and an also integrate seamlessly with the Hadoop Distributed File System. We chose Spark to be the workhorse for the main trading algorithm processing. In order to leverage the features offered by Spark, we used pySpark, which is the Spark API in the Python scripting language. Our algorithm processing engine is responsible for detecting all the trading symbols that are involved in a user-defined trading strategy, and to retrieve the corresponding data from the Hadoop data store. After retrieving the required time series, the engine then runs the rules specified in the strategy against these series and computes the basic returns for the overall portfolio. Once the returns have been computed for the entire data, further portfolio performance metrics are computed. These include the mean return, standard deviation, maximum drawdown and Sharpe Ratio. The final returns time series

and the metrics are then written to temporary output files on disk for use further downstream in the processing pipeline.

### D. Javascript Server Back End

We sought to build a simple web app which could concisely present the results of algorithm backtesting to the user. In order to do this, we implemented a very basic server using the node.js JavaScript runtime.

### E. Visualization

In order to visualize the performance of user defined trading algorithms viz. the benchmark, we built a JavaScript front end for our web application. In order to produce detailed and informative plots, we utilized the Highcharts and Highstock charting libraries.
The schematic in Figure 2. shows an outline of the system design.

## IV. Algorithm (you can change this section name)

Show algorithm and tools you have used to solve your problem.

## V. Software Package Description

Describe the software package that is going to be open sourced. Show some screen shots of how user use it and or UI.

## VI. Experiment Results

Describe the experiment results of your algorithm. Show how did you evaluate the performance of your algorithm.

## VII. Conclusion

Drew a conclusion of your project, describe the contributions of each team member in percentage and discuss future works.

## Acknowledgment

The authors would like to thank... more thanks here

Figure 1. HDFS directory listing showing the flat directory structure for the symbol price data from the QuantQuote dataset.
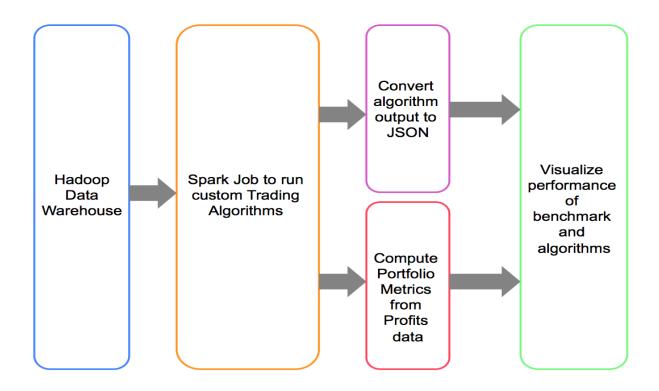


Figure 2. Schematic depicting the overall system organization

APPENDIX

REFERENCES

[1] H. Kopka and P. W. Daly, *A Guide to LaTeX*, 3rd ed. Harlow, England: Addison-Wesley, 1999.