



# A Study in Hadoop Streaming with Matlab for NMR data processing

Kalpa Gunaratna<sup>1</sup>, Paul Anderson<sup>2</sup>, Ajith Ranabahu<sup>1</sup>  
&  
Amit Sheth<sup>1</sup>

<sup>1</sup> **Kno.e.sis** - Ohio Center of Excellence in Knowledge-Enabled Computing  
Wright State University, Dayton, Ohio

<sup>2</sup> **Air Force Research Laboratory**, Biosciences & Protection Division  
Wright-Patterson AFB, Dayton, Ohio



# Outline

- Introduction
- Background
- Design
- Implementation
  - Baseline correction
  - Hadoop streaming
- Results & Discussion
- Conclusion



Kno.E.SIS

# Introduction



FROM INFORMATION TO MEANING

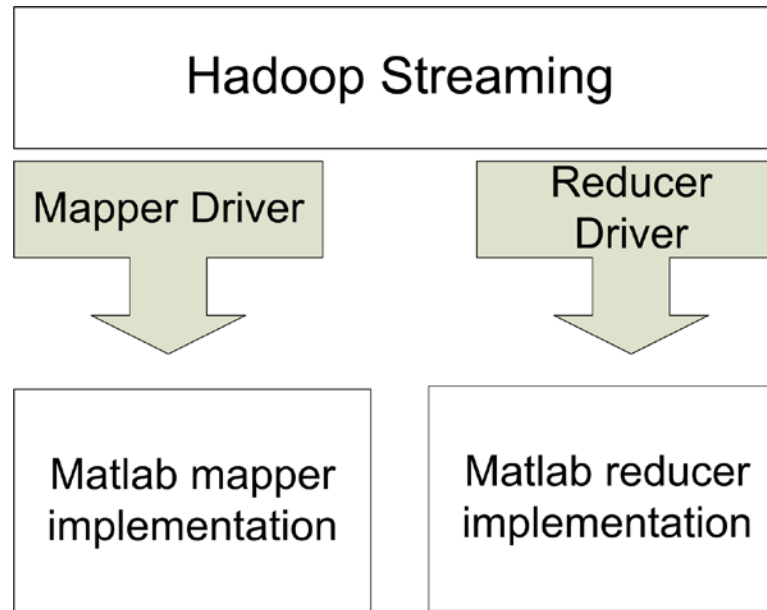
- Biologists confront with huge amount of data (NMR spectrometers, etc).
- Have to undergo numerical processing like baseline correction, normalization, etc. even before doing anything useful.
- Important observation in a Biologists' context,
  - Even though increase in distributed computing tools they avoid using them much.
  - User friendly and domain specific tools are preferred over their lack of performance.
- *Best of both worlds for biologists....*
- Matlab code is run on Hadoop.

- NMR (Nuclear Magnetic Resonance) data analysis normally consists of Giga bytes of data files.
  - A typical  $^1\text{H}$  NMR or  $^{13}\text{C}$  spectrum contain thousands of resonances.
- Metabolomics
  - Assess end product unlike proteomics and genomics.
  - NMR spectroscopy of biofluids is an effective method for identifying variations in states.



# Background cont.

- Baseline distortion
  - Arise from hardware and processing sources.
  - Can lead to incorrect metabolites quantification which leads to spurious scientific conclusions.



- Hadoop streaming is used with C++ driver applications.



kno.e.sis

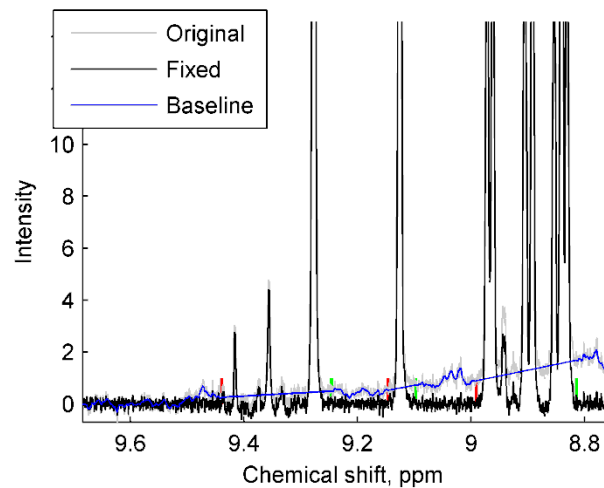
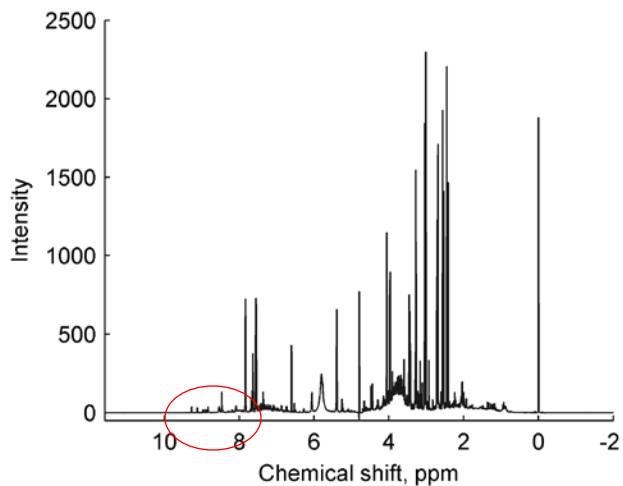
# Design cont.



FROM INFORMATION TO MEANING

- Driver applications are used to read data from the source and call Matlab functions.
- The driver application is responsible for calling relevant Matlab code segments for computations.

- Baseline correction







kno.e.sis

# Implementation cont.

WRIGHT STATE  
UNIVERSITY

FROM INFORMATION TO MEANING

- Baseline correction
  - Whittaker Smoother algorithm is used.
  - The algorithm is written completely in Matlab.

- NMR Data Streaming
  - Driver application is written in C++.
  - Matlab code is compiled with C++ to create a shared library.
  - Driver acts as an interface for mapper in Hadoop and calls Matlab function.
- NMR spectra consist of columns and hence it is inverted to a row oriented file (Hadoop reads line by line).
- Our original Matlab baseline correction desktop code version is trivially changed here.



- Driver creates a relevant Matlab object for a column and passes to the Matlab function.
- For this specific example, a reducer is not necessary since each spectrum is restricted to a single row.
  - If spread across rows, reducers may be needed to format the output.



Kno.E.SIS

# Implementation cont.



FROM INFORMATION TO MEANING

- Technical issues
  - Matlab seemed to have problems with reading directly from Hadoop streaming (need of driver application).
  - Matlab instances need to be available in nodes.



# Results & Discussion

Size	Single machine(sec)	Cluster (sec)
292 KB (1 spectrum)	22	46
2.9 MB (10 spectra)	192	152
28.6 MB (100 spectra)	1996	1563
42.9 MB (150 spectra)	3059	2100
57.2 MB (200 spectra)	4027	2780

Cluster – 16 nodes of Quad core AMD Opteron with 16 GB of RAM

Single machine – 3 GHz dual core CPU with 4 GB of RAM



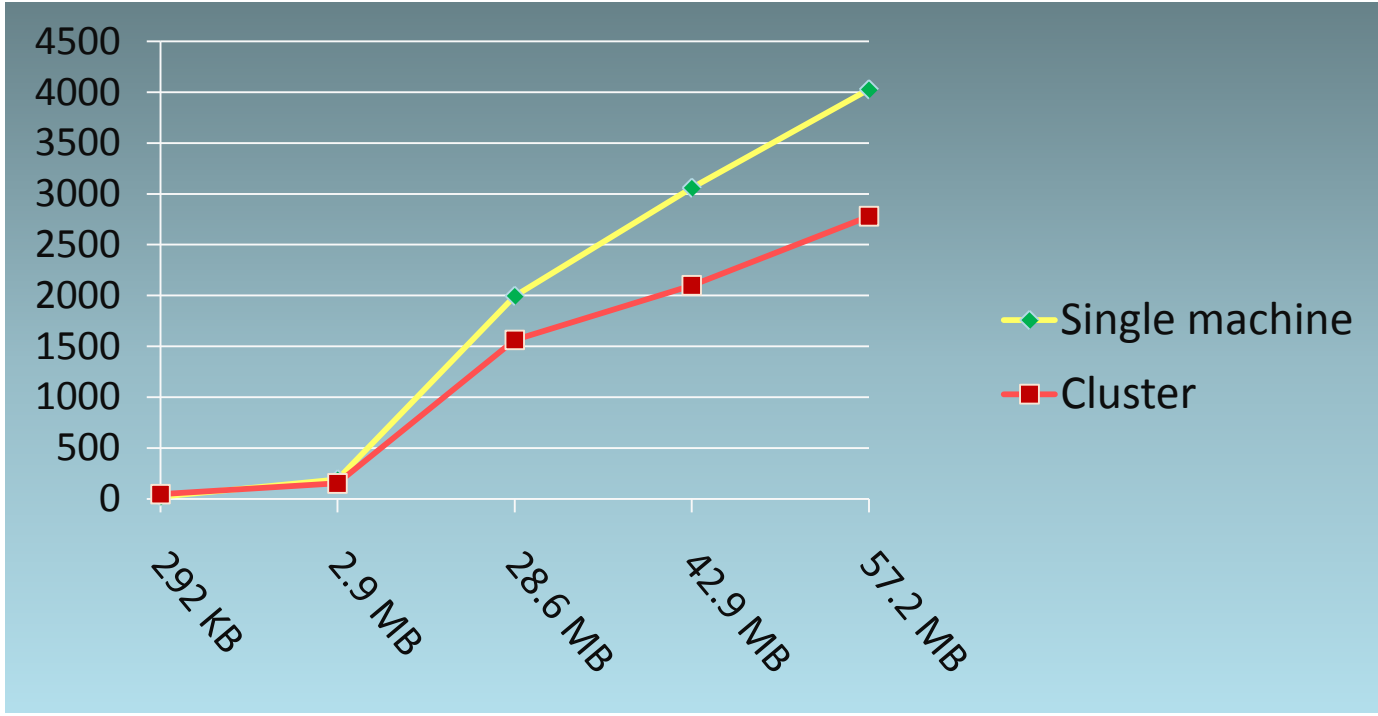
kno.e.sis

Ohio Centers of Excellence

# Results cont.



FROM INFORMATION TO MEANING



- Advantages of using Matlab on Hadoop.
  1. Scientists are relieved from learning new technologies having sharp learning curves (sometimes scripting languages are even incompatible with requirements of biologists).
  2. Non distributed code implementations which are readily available could be used in cloud environment without significant change.
- No need of paradigm shift. Code adoption cost is often expensive and repetitive.
- Facilitates **rapid testing** and **prototyping** where necessary.



# Conclusion

- Cloud computing would not be feasible for scientists if they have to deviate from their routine practices significantly.
- Hence Hadoop streaming allows to use existing Matlab programs in Hadoop clusters.
- Our experiment reflects that using Matlab in Hadoop is feasible and could be extended for various requirements.





Kno.E.SIS

# Questions

WRIGHT STATE  
UNIVERSITY

FROM INFORMATION TO MEANING





kno.e.sis

Ohio Centers of Excellence

WRIGHT STATE UNIVERSITY

FROM INFORMATION TO MEANING

Thank You!

<http://knoesis.org>

WRIGHT STATE UNIVERSITY

kno.e.sis