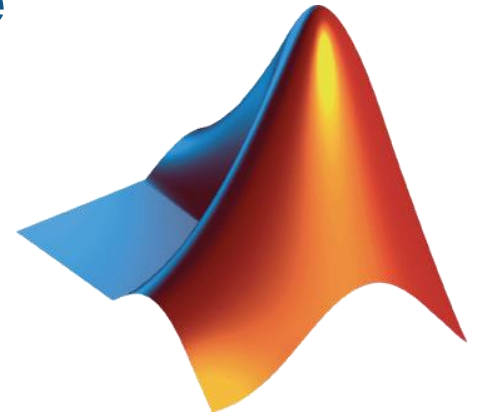# Outils pour l'analyse prédictive parallèle de multiples sources de données non structurées

**Forum Ter@tec**

**Mercredi 25 juin 2015**

**Marc Wolff – Application Engineer HPC & Big Data**

# Challenges Businesses Are Facing Today

**Big Data for evidence-based decision making**



- **Goal**
  - Large (and increasing) amount of available data
  - Leverage data to make better decision

- **Organizational issues**
  - Rapid evolution
  - Data scientist need to share their algorithms and results efficiently

- **Technical issues**
  - Datasets do not fit in the memory of a single computer
  - Processing these data requires huge computing resources

# How MATLAB Helps Tackling Big Data

**Data**

- More data, faster access
- Complex / incomplete / changing formats handling
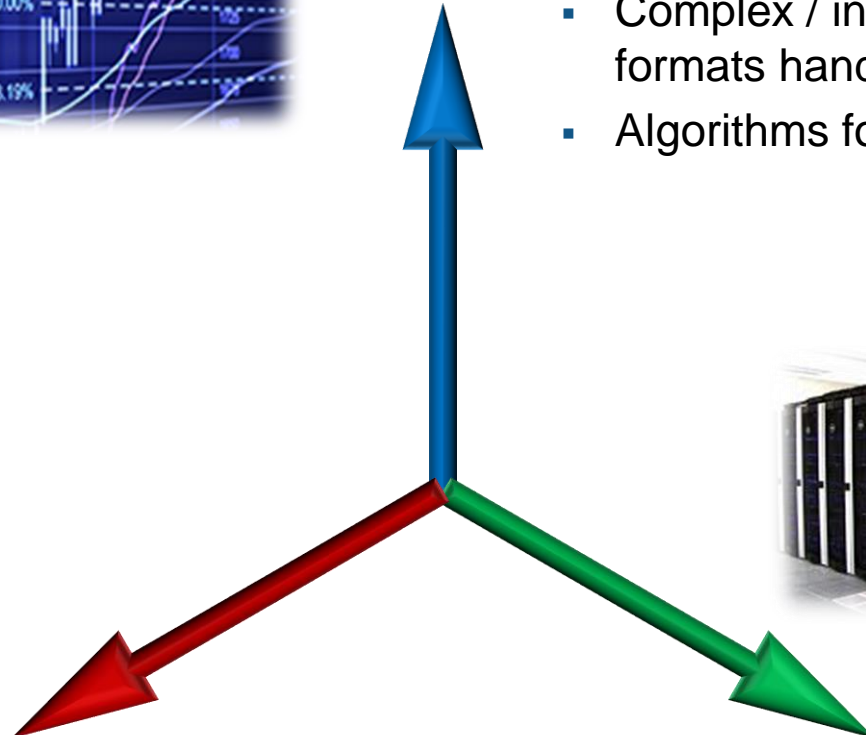- Algorithms for complex problems

**People & Systems**

- Share algorithms & protect IP
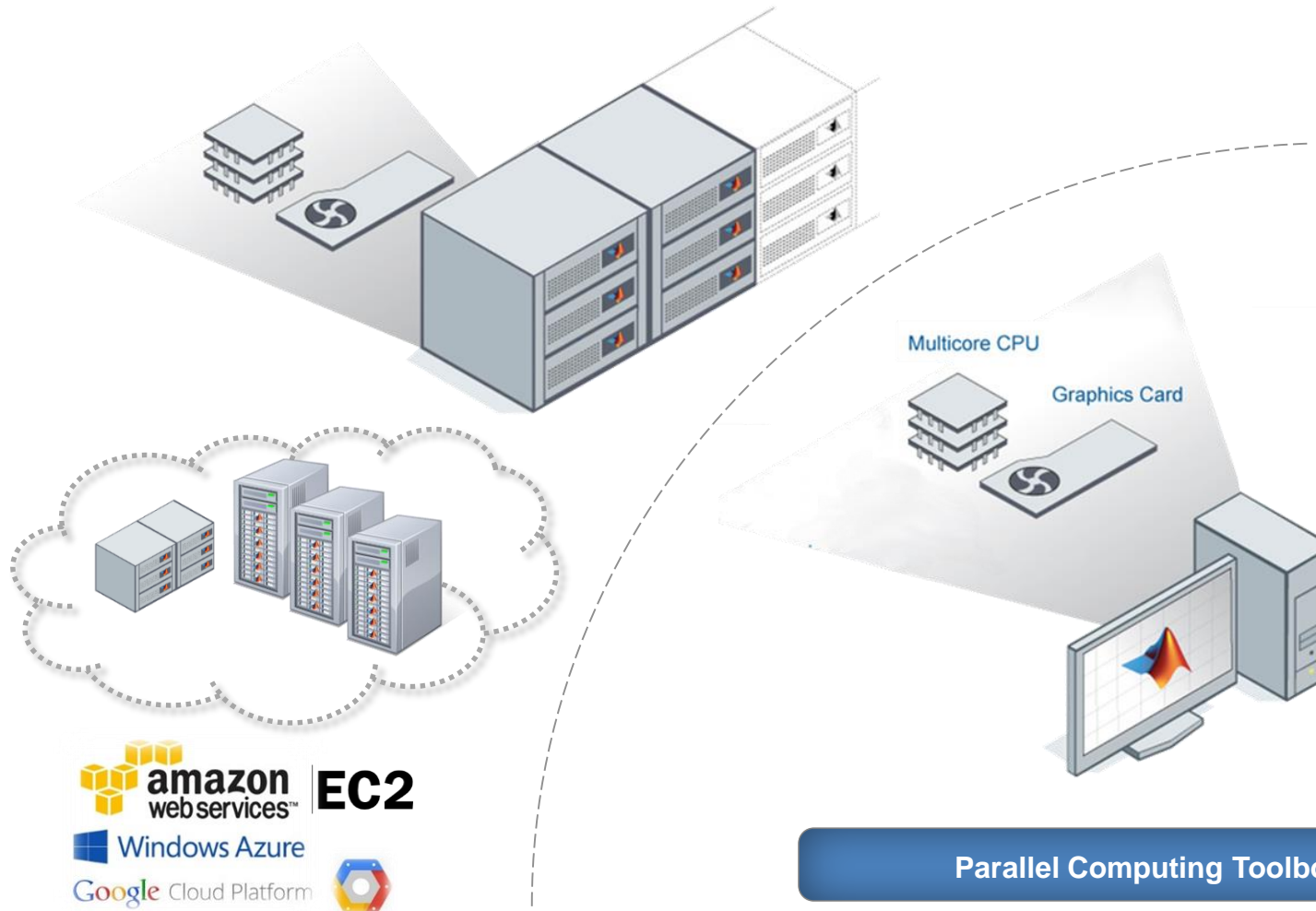- Real-time analytics

**Compute Power**

- Leverage HPC facilities & clouds
- Distributed memory framework

# Scaling Out Calculations

**MATLAB Distributed Computing Server (MDCS)**



Multicore CPU

Graphics Card

amazon web services™ | EC2

Windows Azure

Google Cloud Platform

**Parallel Computing Toolbox (PCT)**

# MATLAB & High Performance Computing?

## Jack Dongarra

From Wikipedia, the free encyclopedia

**Jack J. Dongarra** (born July 18, 1950) is an American University Distinguished Professor of Computer Science in the Electrical Engineering and Computer Science Department[10] at the University of Tennessee. He holds the position of a Distinguished Research Staff member in the Computer Science and Mathematics Division at Oak Ridge National Laboratory, and is an Adjunct Professor in the Computer Science Department at Rice University. Dongarra holds the Turing Fellowship in the schools of Computer Science and Mathematics at the University of Manchester. He is the founding director of Innovative Computing Laboratory.[11][12][1][13][14][15]

**Contents** [hide]

1 Education
2 Research
3 References
4 External links

## Education [edit]

Dongarra received a Bachelor of Science degree in Mathematics from Chicago State University in 1972 and a Master of Science in Computer Science from the Illinois Institute of Technology in 1973. He received his Doctor of Philosophy in Applied Mathematics from the University of New Mexico in 1980 under the supervision of Cleve Moler.[2] He worked at the Argonne National Laboratory until 1989, becoming a senior scientist.

## Research [edit]

He specializes in numerical algorithms in linear algebra, parallel computing, the use of advanced-computer architectures, programming methodology, and tools for parallel computers. His research includes the development, testing and documentation of high quality mathematical software. He has contributed to the design and implementation of the following open source software packages and systems: EISPACK, LINPACK, the BLAS, LAPACK, ScaLAPACK,[3][4] Netlib, PVM, MPI,[5] NetSolve,[6] TOP500, ATLAS,[7] and PAPI.[8] With Eric Grosse, he pioneered the open source distribution of numeric source code via email with netlib. He has published approximately 300 articles, papers, reports and technical memoranda and he is coauthor of several books. He was awarded the IEEE Sid Fernbach Award in 2004 for his contributions in the application of high performance computers using innovative approaches; in 2008 he was the recipient of the first IEEE Medal of Excellence in Scalable Computing; in 2010 he was the first recipient of the SIAM Special Interest Group on Supercomputing's award for Career Achievement; in 2011 he was the recipient of the IEEE IPDPS Charles Babbage Award; and in 2013 he was the recipient of the ACM/IEEE Ken Kennedy Award for his leadership in designing and promoting standards for mathematical software used to solve numerical problems common to high performance computing.He is a Fellow of the AAAS, ACM, SIAM, and the IEEE and a member of the National Academy of Engineering.
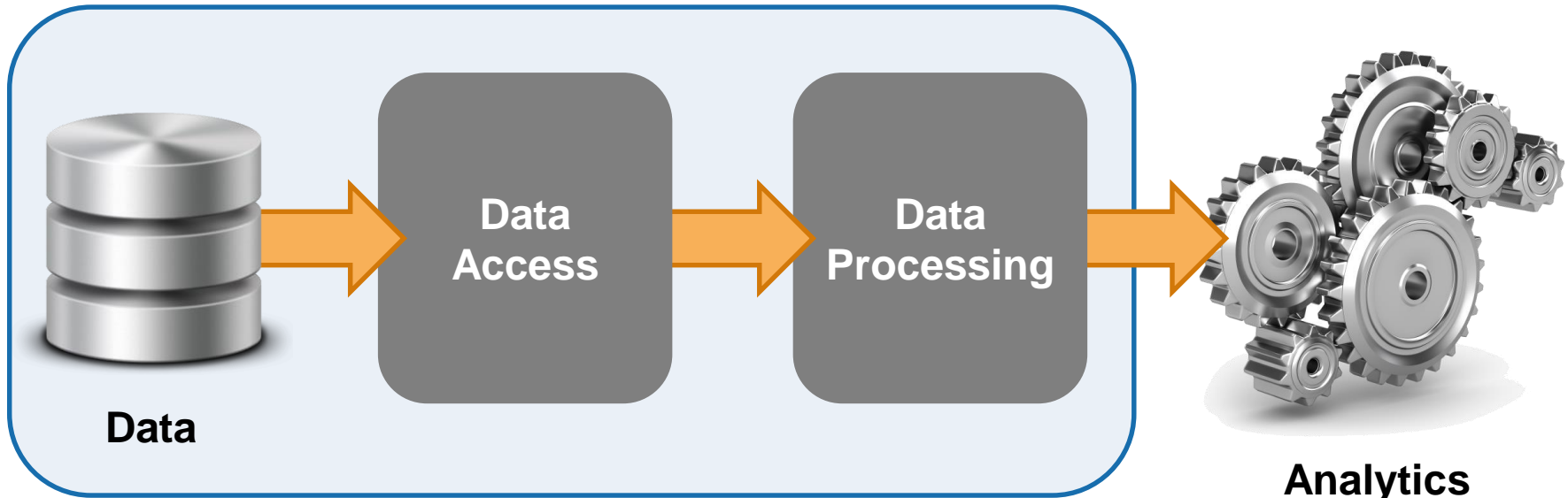
**Jack Dongarra**

Jack Dongarra

| | |
|---|---|
| **Born** | July 18, 1950 (age 64) Chicago |
| **Citizenship** | American / United States |
| **Nationality** | American |
| **Fields** | Computer Science Computational science Parallel computing[1] |
| **Institutions** | University of Tennessee University of New Mexico Argonne National Laboratory Oak Ridge National Laboratory University of Manchester |
| **Alma mater** | University of New Mexico |
| **Thesis** | *Improving the Accuracy of Computed Matrix Eigenvalues* (1980) |
| **Doctoral advisor** | Cleve Moler[2] |

# Data Analytics Workflow



**Data**

**Data Access**
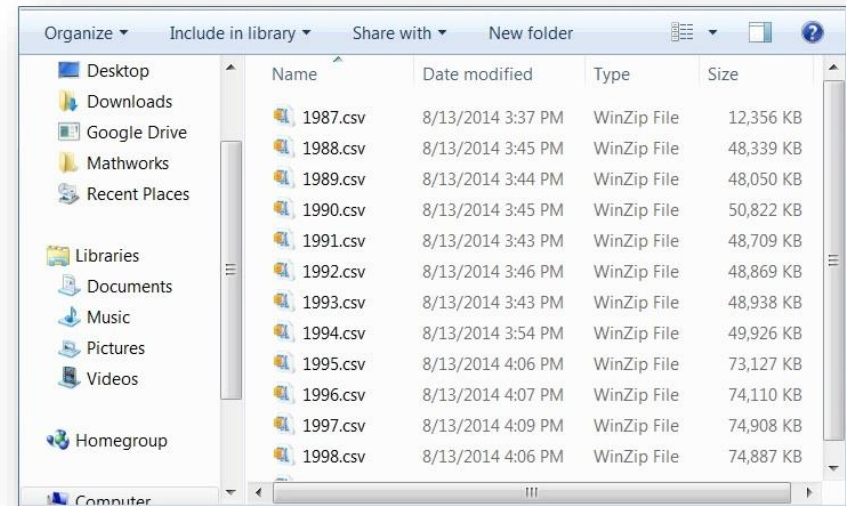
**Data Processing**

**Analytics**

**Today's presentation**

– Descriptive Statistics
– Machine Learning
– Neural Networks
– MapReduce

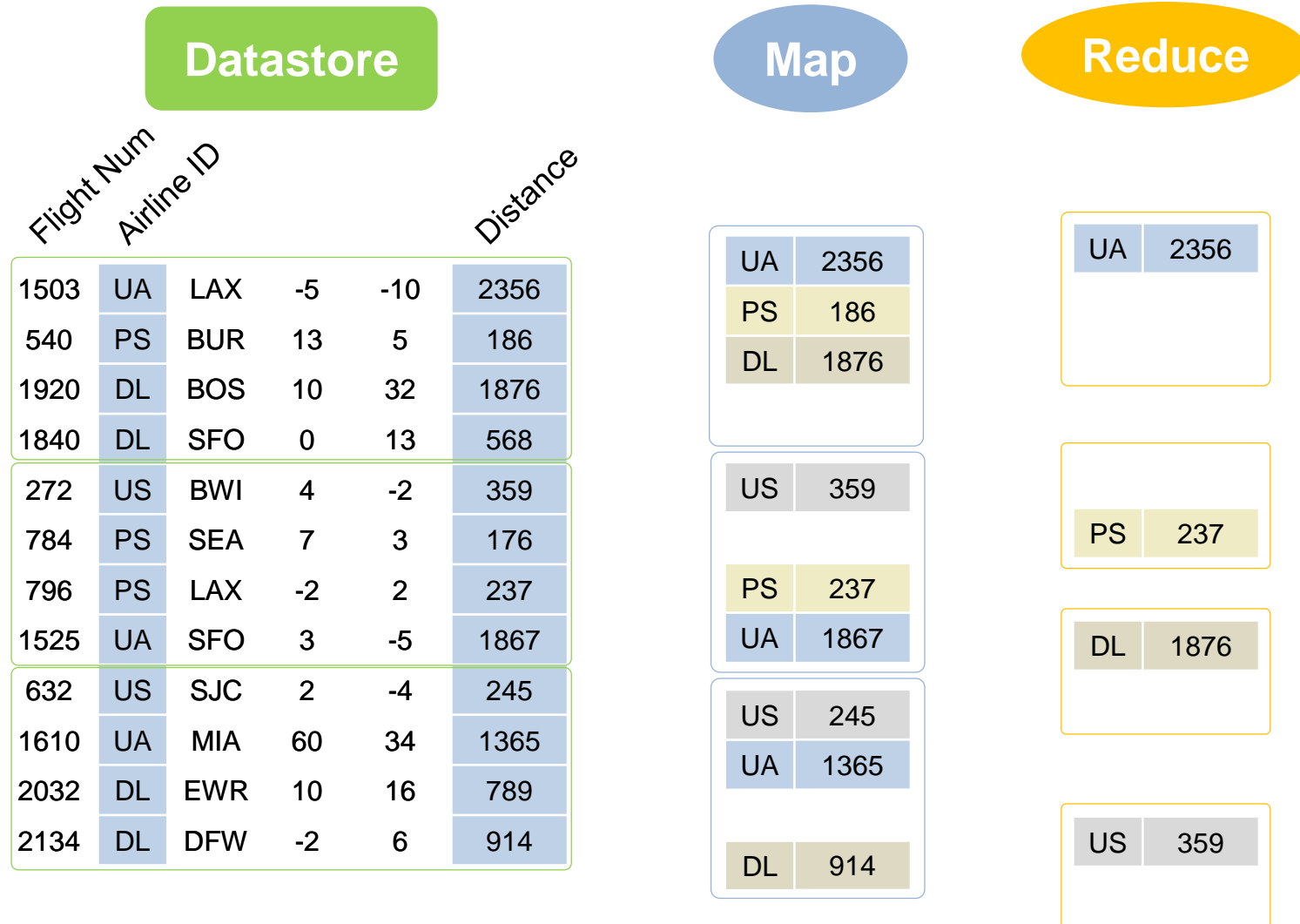# Access Large Datasets
**datastore**

- Data container that allows to easily read data that are **too large to fit in the computer's memory**

- Incremental read: data loaded in memory by parts

- Data sources of various natures
  - Database (using Database Toolbox)
  - Single text file or collection of text files
  - MATLAB is generally able to read and write data directly from/to HDFS

- **datastore** can be partitioned using the **partition** function
  - Allows **parallel** data access
  - Take advantage of parallel file systems

# MapReduce Programming Model
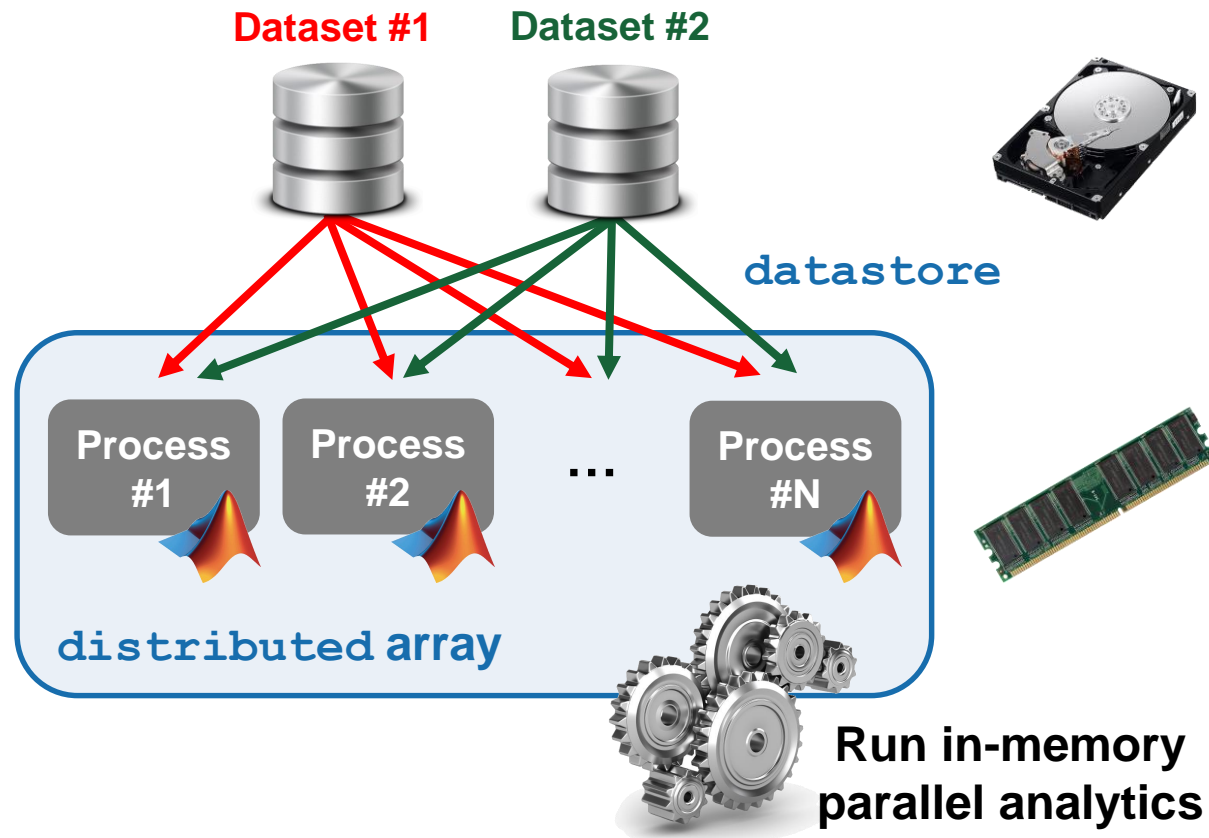**Strengths and limitations**

- `mapreduce` has been introduced in MATLAB **R**2014**b**

- **Strengths**
  – Analytics are made easy when they fit in the MapReduce framework
  – MapReduce on Hadoop can take advantage of data locality in HDFS

**Shared filesystem**

**High-speed network**

- **Limitations**
  – Subset-by-subset data processing, no vision of the whole dataset
  – Scalability issues in some cases

# A Parallel Programming Model for Predictive Analytics

- Parallel computing implementation in MATLAB
  - Capabilities all based on MPI
  - MATLAB offers a transparently distributed data structure: `distributed`



**Dataset #1**    **Dataset #2**

`datastore`

`distributed array`

Process #1   Process #2   ...   Process #N

**Run in-memory parallel analytics**

# A Parallel Programming Model for Predictive Analytics
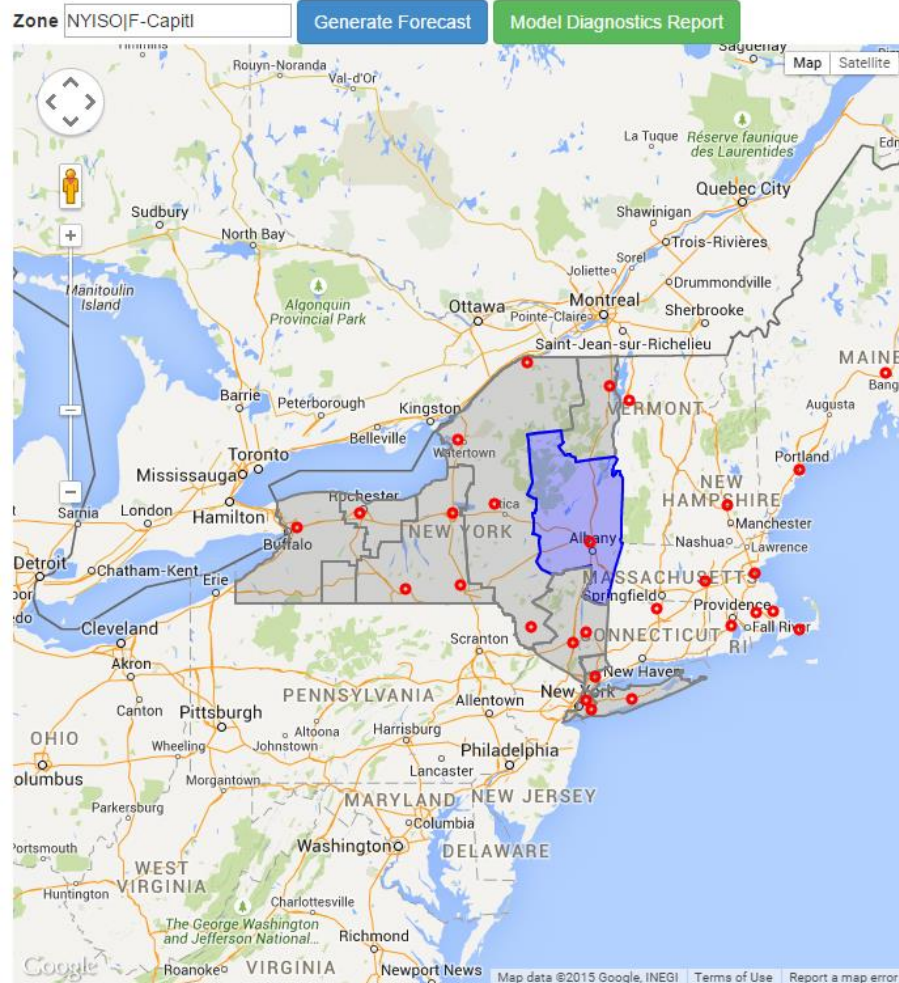## Scalability analysis

| | Does it scale? | Does it require to store the complete dataset on a single process? |
|---|---|---|
| 1. Read data from one or multiple datasets and preprocess it | ✔ | ✔ |
| 2. Store the preprocessed data in a distributed array | ✔ (up to 100s of processes) | ✔ |
| 3. Run in-memory analytics on the distributed array | ✔ (depends on the algorithm) | ✔ (depends on the algorithm) |

# A Parallel Programming Model for Predictive Analytics
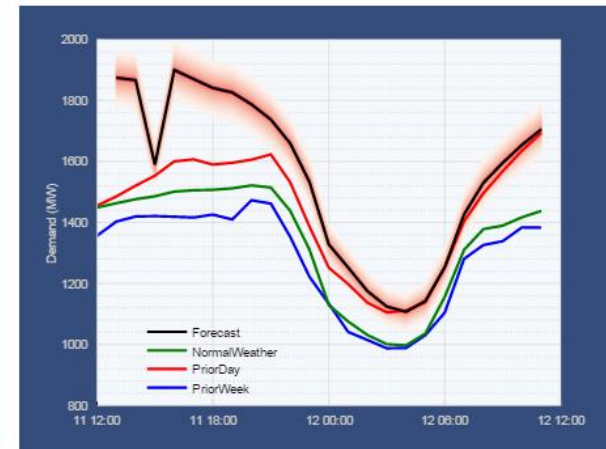## Example: power load forecasting



http://ec2-54-165-201-58.compute-1.amazonaws.com:8080/DemandForecastWeb/

# A Parallel Programming Model for Predictive Analytics
## Example: power load forecasting

- **Goal**
  - Develop a predictive model to forecast electrical power consumption
  - Deploy the prediction tool in power plants to adjust production

- **Predictors from different datasets**
  - Power consumption over the previous days
  - Calendar information: day of the week? holiday period?
  - Climate data

- **Challenges**
  - Multiple data sources with different formatting
  - Datasets have different samplings

- **Final result**
  - Predictive model based on Neural Networks
  - Deployed in production using MATLAB Production Server

# Key Takeaways

- Access large datasets with MATLAB
  - `datastore` allows to read datasets that do not fit in memory
  - `partition` allows parallel data access

- Tools for easily developing algorithms and scaling out computations
  - No need to be an expert in computer science
  - No need to be an expert in parallel computing

- MathWorks development teams heavily involved
  - Continuous development and improvement
  - New features in each release

## Thank you!

## Q & A session