



The Big Picture of Data Science Approaches and Tools

AMARNATH GUPTA

Three Definitions

Data Intensive Computing

a class of parallel computing applications which use a data parallel approach to processing large volumes of data

Big Data

data sets so large, diverse, complex or fast that traditional ingestion strategies, data processing and computing algorithms, and system infrastructure are inadequate

Data Science

almost everything that has something to do with data: collecting, analyzing, modeling to model, predict and infer information for domain applications

Data Science: an Action Plan for Expanding the Technical Areas of the Field of Statistics

William S. Cleveland

Statistics Research, Bell Laboratories, 600 Mountain Avenue, Murray Hill NJ07974, USA
E-mail: wsc@research.bell-labs.com

Summary

An action plan to enlarge the technical areas of statistics focuses on the data analyst. The plan sets out six technical areas of work for a university department, and advocates a specific allocation of resources devoted to research in each area and to courses in each area. The value of technical work is judged by the extent to which it benefits the data analyst, either directly or indirectly. The plan is also applicable to government research labs and corporate research organizations.

Key words: Future; Applications; Computing; Methods; Models; Theory.

1 Summary of the Plan

This document describes a plan to enlarge the major areas of technical work of the field of statistics. Because the plan is ambitious and implies substantial change, the altered field will be called “data science”.

The focus of the plan is the practicing data analyst. A basic premise is that technical areas of data science should be judged by the extent to which they enable the analyst to learn from data. The benefit of an area can be direct or indirect. Tools that are used by the data analyst are of direct benefit. Theories that serve as a basis for developing tools are of indirect benefit. A broad successful theory can have a wide-ranging benefit, affecting data analysis in a fundamental way. For example, the Bayesian theory of statistics affects all methods of estimation and distribution.

The plan sets out six technical areas for a university department, and advocates a specific allocation of resources to research and development in each area as a percent of the total resources that are

The “official” start of Data Science

(25%) **Multidisciplinary Investigations:** data analysis collaborations in a collection of subject matter areas.

(20%) **Models and Methods for Data:** statistical models; methods of model building; methods of estimation and distribution based on probabilistic inference.

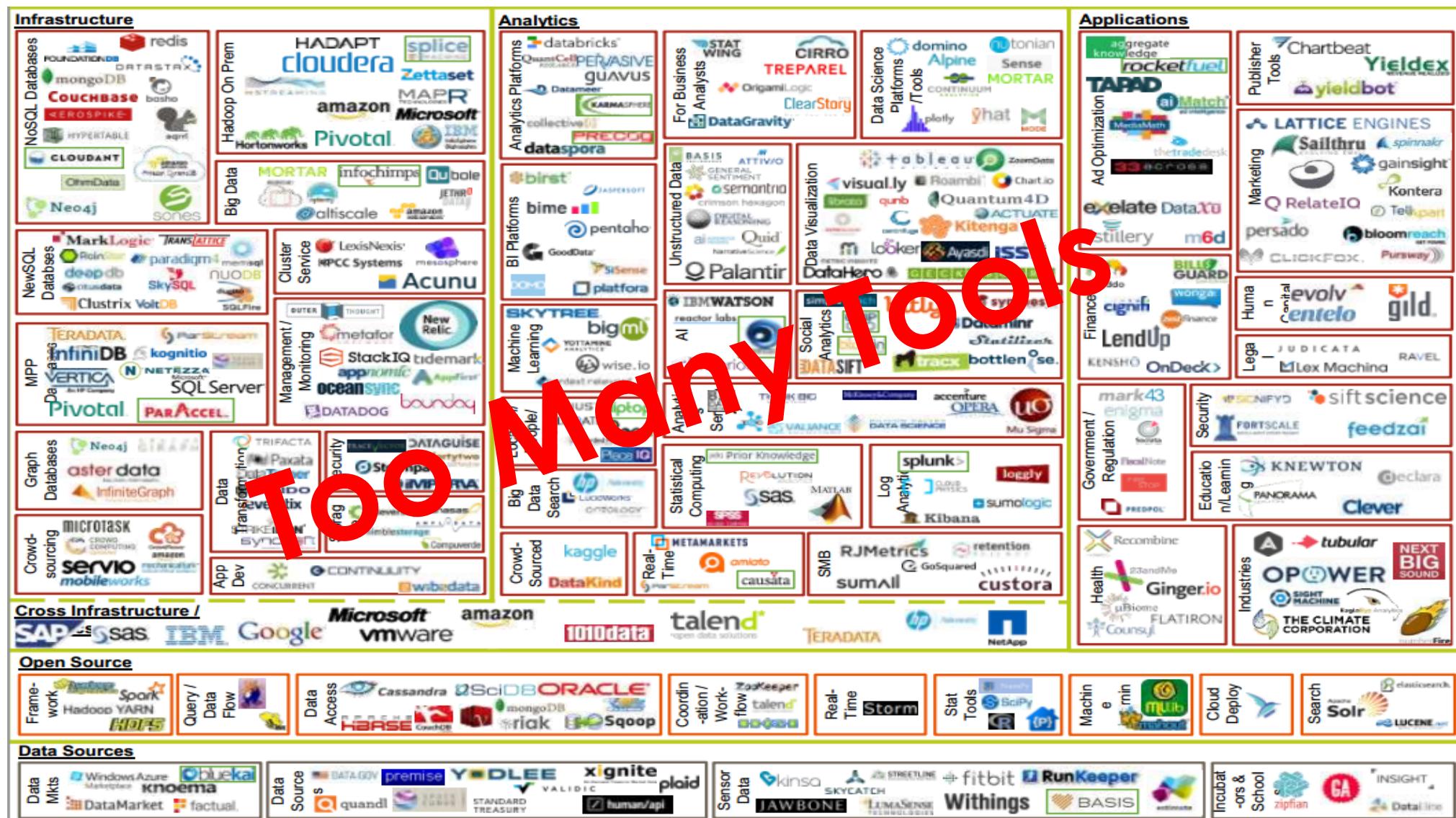
(15%) **Computing with Data:** hardware systems; software systems; computational algorithms.

(15%) **Pedagogy:** curriculum planning and approaches to teaching for elementary school, secondary school, college, graduate school, continuing education, and corporate training.

(5%) **Tool Evaluation:** surveys of tools in use in practice, surveys of perceived needs for new tools, and studies of the processes for developing new tools.

(20%) **Theory:** foundations of data science; general approaches to models and methods, computing with data, teaching, and tool evaluation; mathematical investigations of models and methods, computing with data, teaching, and evaluation.

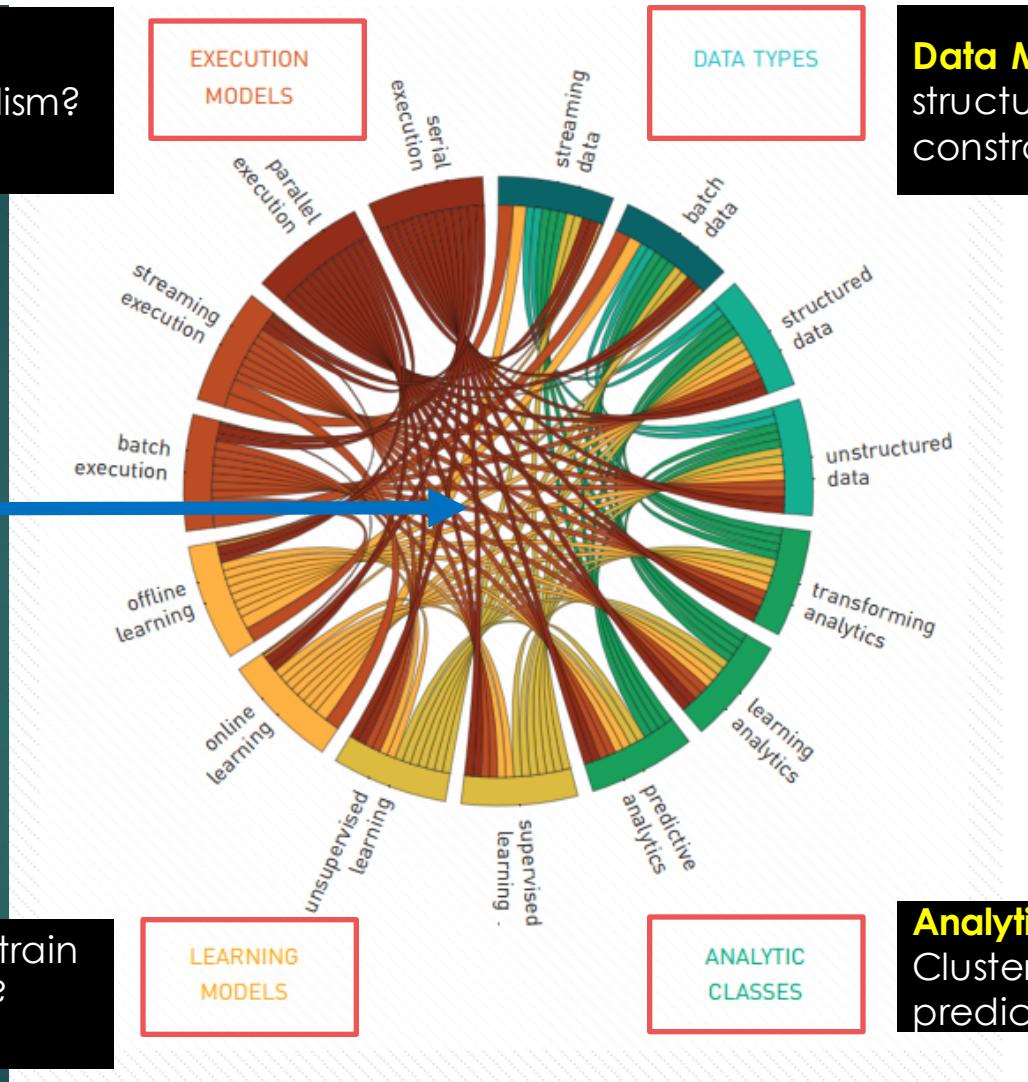
Domain Knowledge + Analytical Methods + Algorithms+ Databases



Execution Model: Data parallelism? Task parallelism? Stream processing?

There are tools for each of these combinations

Learning Model: How to train data? Model of training? Algorithm of learning?



Data Model: specifies the structure, operations and constraints on data

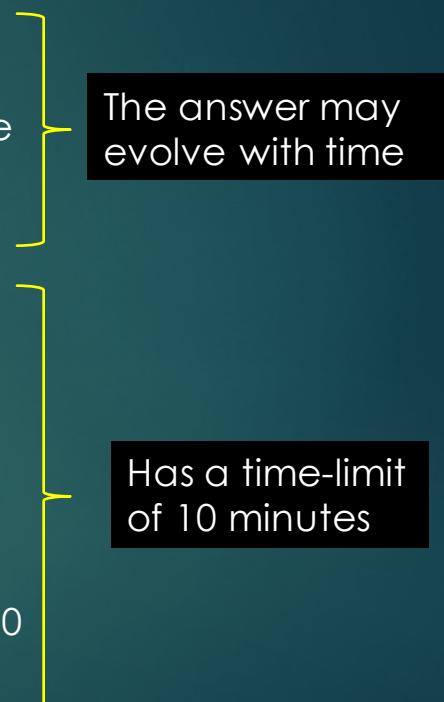
Partitioning the Problem

Analytic Class: Classification? Clustering? Time-series based prediction? Recommendation?

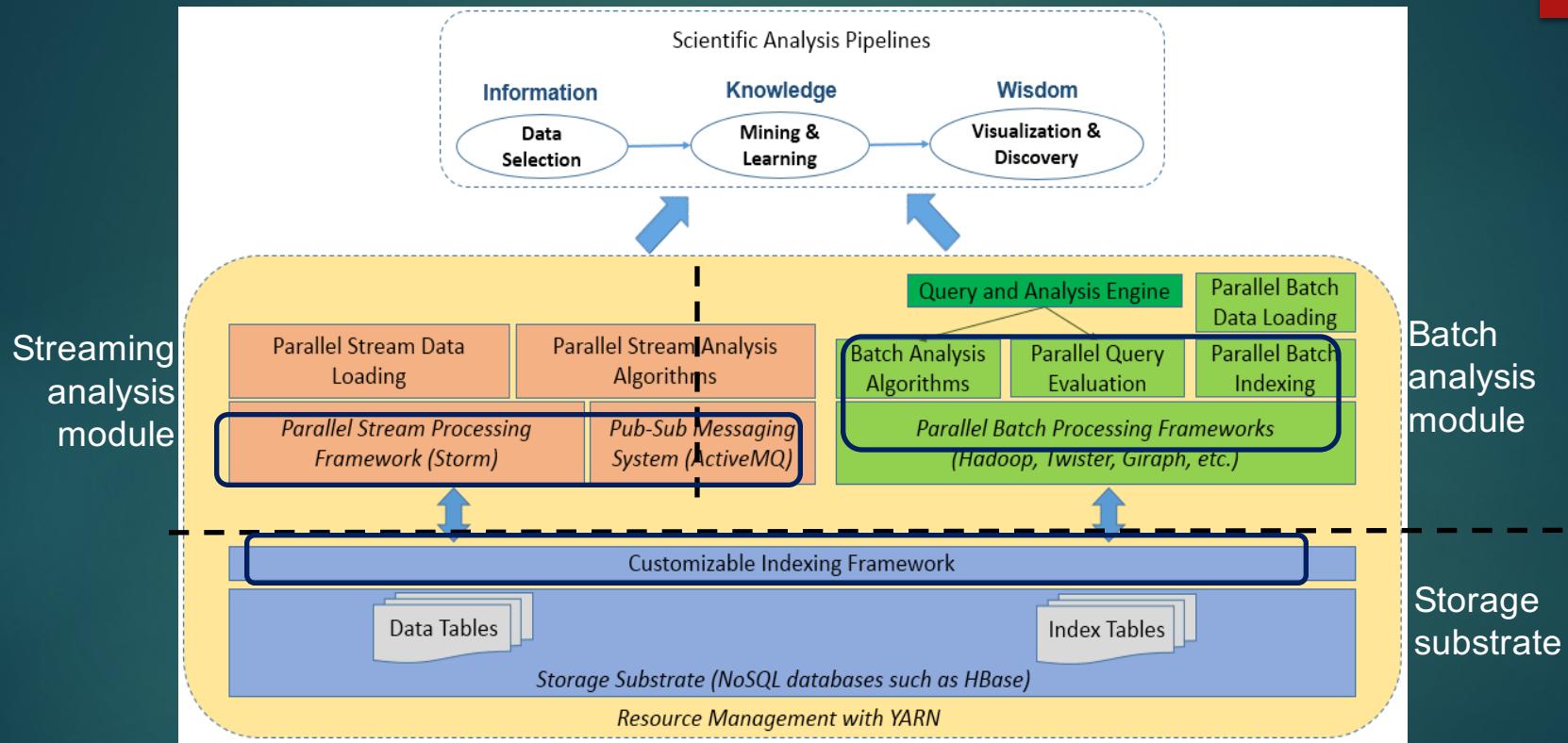
Interconnection Among the Component Parts of Data Science

Source: Booz Allen Hamilton

Some Data Science Problems

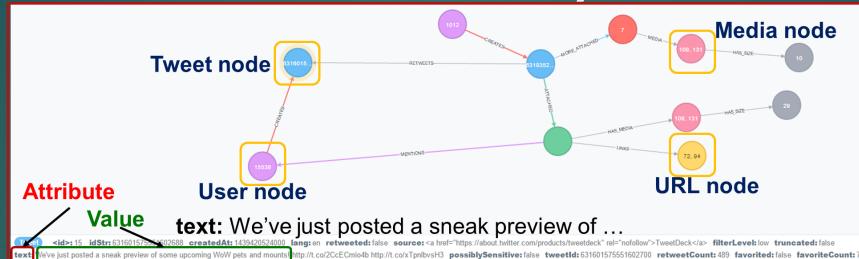
- ▶ Looking at public tweets, is it possible to
 - ▶ Predict how the political leanings of the tweeting populations are changing?
 - ▶ Predict the success or failure of a political campaign
 - ▶ Given
 - ▶ Historical information
 - ▶ Current weather data
 - ▶ Current state of wildfires
 - ▶ Simulation engines for wind patterns and fire perimeters
 - ▶ Can we predict if the fire is likely to reach *this area* in the next 10 minutes?
- 
- The answer may evolve with time
- Has a time-limit of 10 minutes

A General Architecture



Ideally, must run on a (virtual) cluster of machines to offer horizontal scalability

Data Genres (Data following different data models) - I



- ▶ A tweet is a complex data object
 - ▶ Data format – JSON MongoDB?
 - ▶ Metadata – relational (i.e., table like) Any relational DBMS that scales?
 - ▶ Tweet body – text Index with Apache Solr?
 - ▶ Inter-tweet connections – network Neo4J?
 - ▶ Inter-user connections – network Spark/GraphX?
- ▶ Number of tweets we collect – 4-8GB/day, >1 TB over 5 months

} Any Data Mgmt System that works well on all data genres?

Not really!!

Analytical Query Processing

- ▶ Do we need
 - ▶ More custom data processing, or Hadoop or Spark-based Solutions
 - ▶ Store data and perform arbitrary queries, or Vertical (relational), AsterixDB (JSON, XML), VoltDB (in-memory)...
 - ▶ Both at the same time SPARK-SQL ??
 - ▶ For our data science example
 - ▶ Store in AsterixDB
 - ▶ Network queries in Neo4J
 - ▶ Compute in SPARK
 - ▶ Move data with a message bus (e.g., RabbitMQ, Apache Ignite, ...)
- 
- Multi-System
Solution

Analytics Models

- ▶ ~ 200 Million Tweets
- ▶ We need
 - ▶ Classifiers to identify only the political tweets
 - ▶ Train to eliminate ads, personal tweets, spams
 - ▶ Size of training set – 20k
 - ▶ Topic Model construction over time windows
 - ▶ Vocabulary size – 30,000
 - ▶ Document count – 170 million
 - ▶ Multivariable Regression

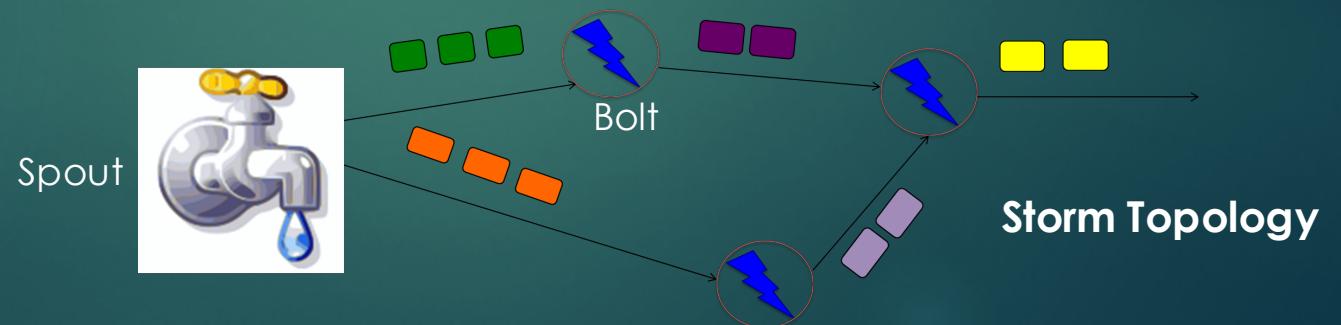
Spark with MLLib running over a virtual cluster of machines

SparkR to use the R statistical analysis system over Spark

Why Spark?
Iterative Memory-intensive Data Parallel Computation

Execution Model – Stream Processing

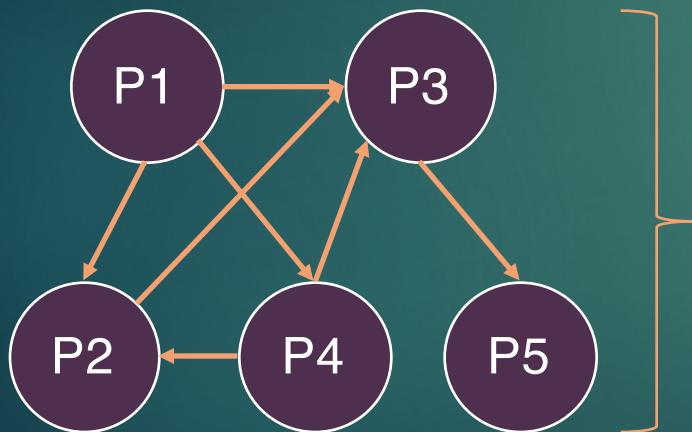
- ▶ Apache Storm
 - ▶ Operations that can be performed
 - ▶ **Filter:** forward only tuples which satisfy a condition
 - ▶ **Joins:** When receiving two streams A and B, output all pairs (A,B) which satisfy a condition
 - ▶ **Apply/transform:** Modify each tuple according to a function



Execution Model – Network Processing with Bulk Synchronous Parallelism



- Processors
 - Have local memory
 - Can perform some computation



- Processors can communicate pairwise
- Communication can overlap with another node's computation

Barrier Synchronization

Apache Giraph

Spark/GraphX

Conclusion

- ▶ No universal big picture solution for large-scale data science problem
- ▶ Tool choice based on data types, analytical model, execution model and learning model
- ▶ A glimpse of tools to consider within the context of one data science problem