

# The Big Data team

Marcus Berg



AMERICAN ASSOCIATION FOR PUBLIC OPINION RESEARCH

## **AAPOR Report on Big Data**

**AAPOR Big Data Task Force**

*February 12, 2015*

# **Task Force Members:**

*Lilli Japac, Co-Chair, Statistics Sweden*

*Frauke Kreuter, Co-Chair, JPSM at the U. of Maryland, U. of Mannheim & IAB*

*Marcus Berg, Stockholm University*

*Paul Biemer, RTI International*

*Paul Decker, Mathematica Policy Research*

*Cliff Lampe, School of Information at the University of Michigan*

*Julia Lane, American Institutes for Research*

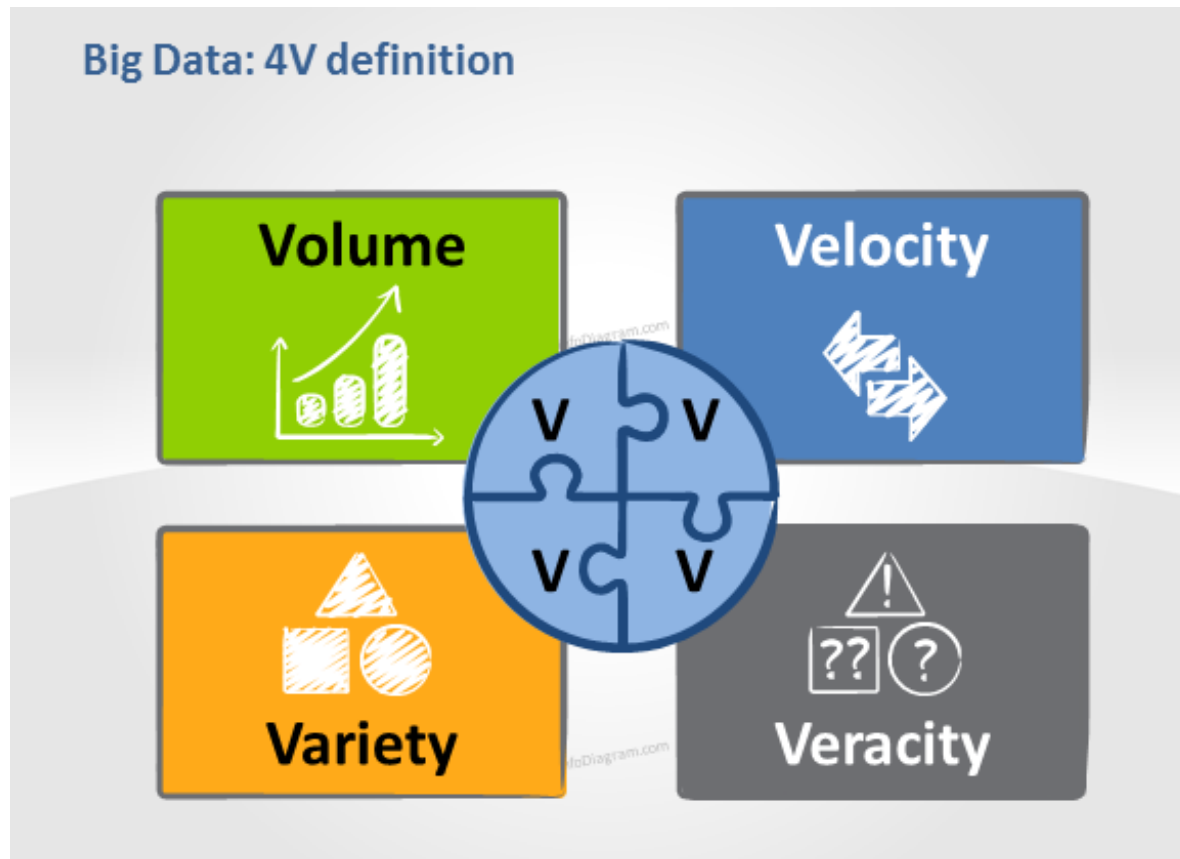
*Cathy O'Neil, Johnson Research Labs*

*Abe Usher, HumanGeo Group*

**Acknowledgement: We are grateful for comments, feedback and editorial help from Eran Ben-Porath, Jason McMillan, and the AAPOR council members.**

# Another way to visualize it

*because visualization is important*

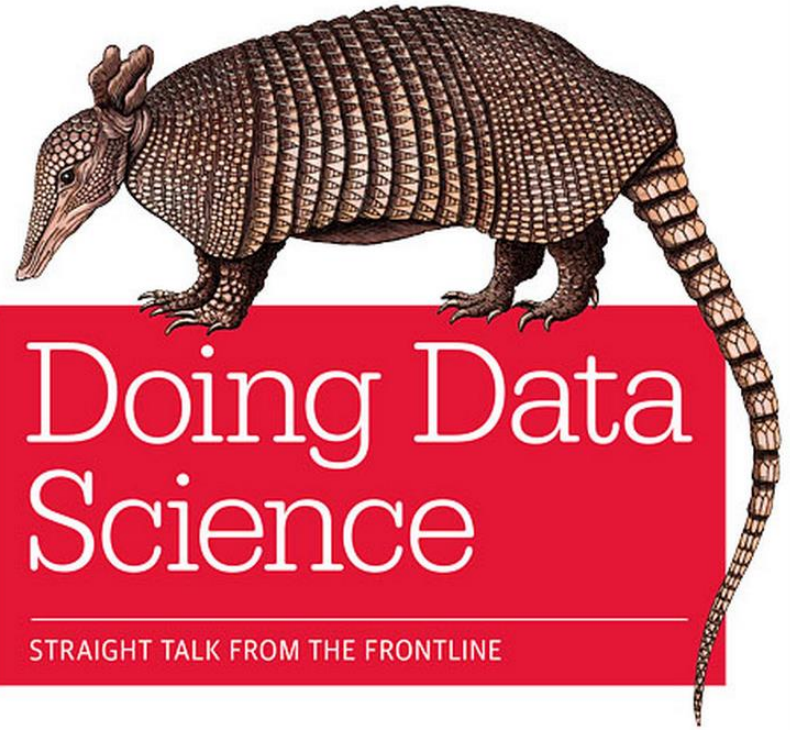


# Big Data

Data science

Statistics

O'REILLY®



Doing Data  
Science

STRAIGHT TALK FROM THE FRONTLINE

Rachel Schutt & Cathy O'Neil

# Four roles of a Big Data team

*each provide different skills*

## DOMAIN EXPERT

User, analyst, or leaders with deep subject matter expertise related to the data, its appropriate use, and its limitations

## SYS ADMIN

Team member responsible for defining and maintaining a computation infrastructure that enables large scale computation



## RESEARCHER

Team member with experience applying formal research methods, including survey methodology and statistics

## COMPUTER SCIENTIST

Technically skilled team member with education in computer programming and data processing technology

# Tech roles

*Computer  
Scientist*

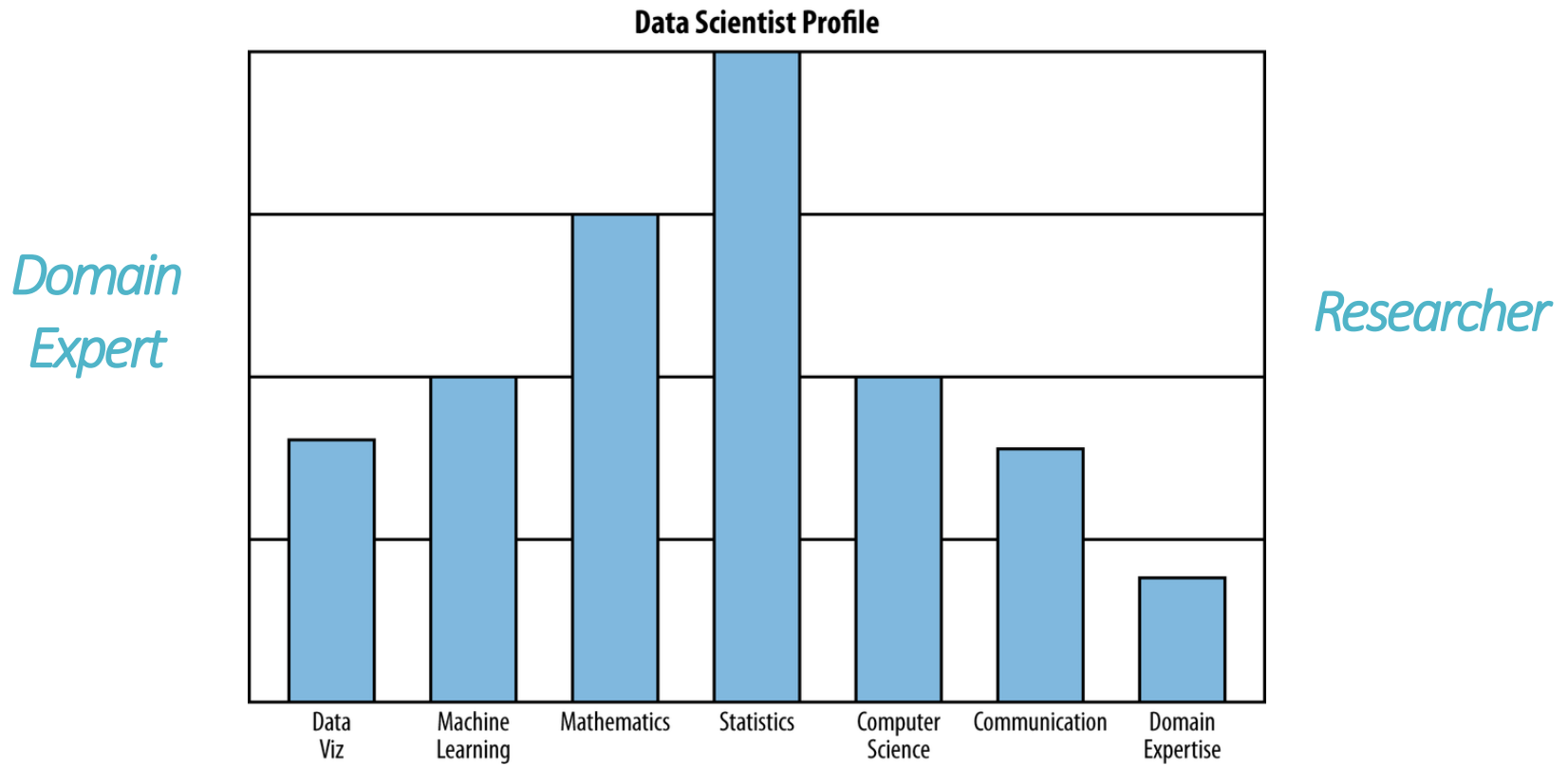
Distributed File Systems (DFS)  
*A 'doubly secure' storage solution*

*System  
Administrator*

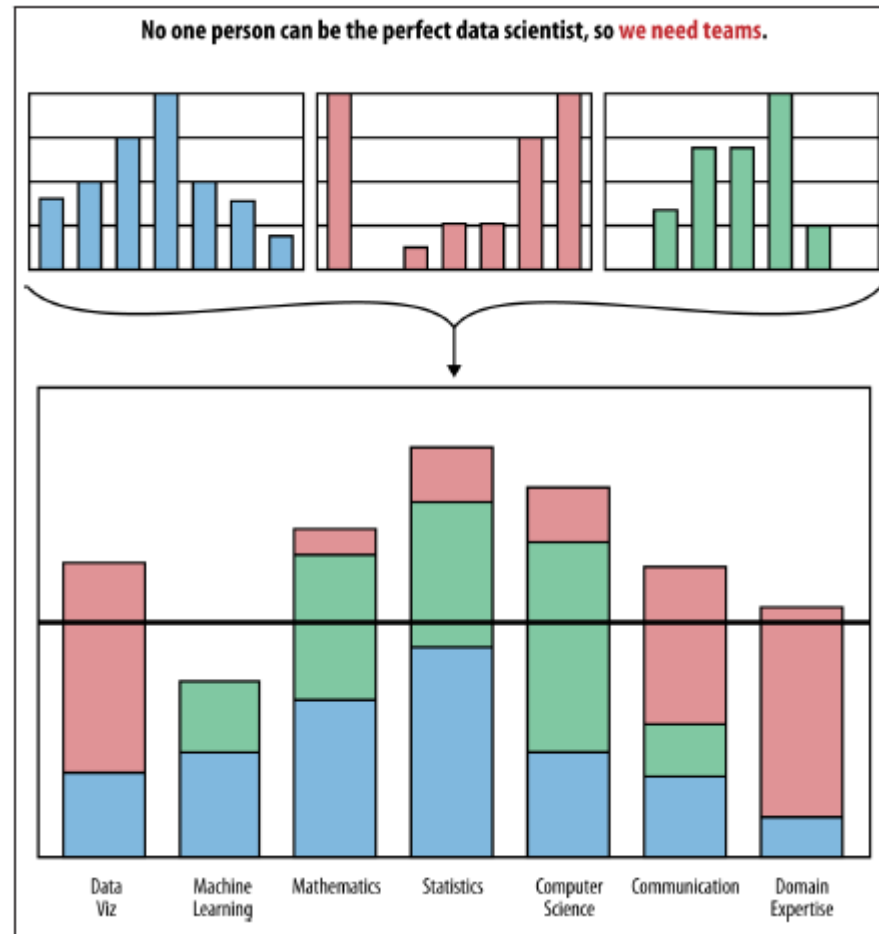
MapReduce (MR)  
*A 'divide and conquer' data processing model*



# Non-tech roles



# The data scientist team





# Big Data: Layers

---

Data Output

Example: map visualization

Data Analysis

Example: Hadoop MapReduce

Data Storage

Example: Hadoop Distributed File System

Data Source(s)

Examples: geolocated social media  
(Proxy variable for behavior of interest)

Thank you!

[malex.berg@gmail.com](mailto:malex.berg@gmail.com)