

DATA SCIENCE ECOSYSTEM

M. TAMER ÖZSU

U. WATERLOO

NANCY REID

U. TORONTO

RAYMOND NG

UBC

DATA SCIENCE/BIG DATA IN THE NEWS...

CBC | MENU ▾

Big Brother meets Big Data, in an office near you

The Atlantic

Sponsor Content:
What's this?

THE WALL STREET JOURNAL.

CIO JOURNAL
Carnival Strategy Chief Bets That Big Data Will Optimize Prices

New York Times Adapts Data Science Tools for Advertisers

Team will help lure marketers with tools to predict which articles will resonate with certain readers to better target advertising

Forbes / Tech

MAY 27, 2015 @ 10:20 AM

34,550

How Big Data And The Internet Of Things Improve Public Transport In London

The Little Black Book of Billionaires

Improve Public

npr

SCIENCE

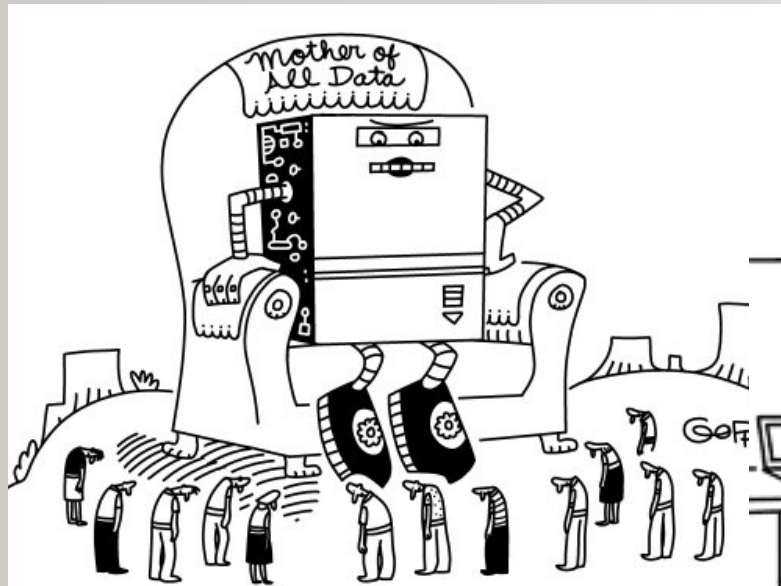
The Big Idea Behind Big Data

BIG DATA AND HOLLYWOOD: A LOVE STORY

Data Veracity is Critical for Insurers to Make Better Business Decisions, According to Accenture Report

Français

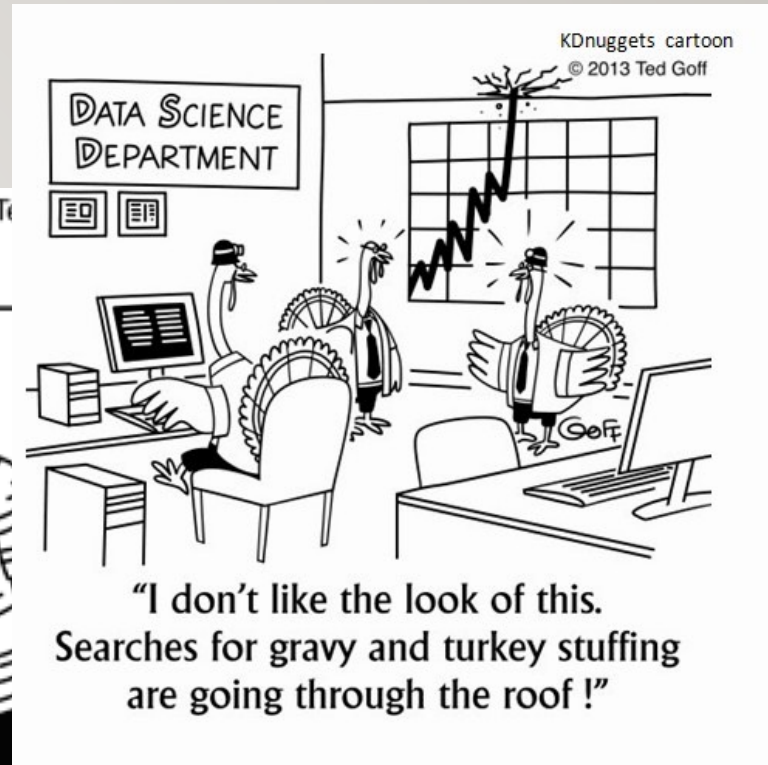
DATA SCIENCE EVERYWHERE!...



"No candy? No flowers? No cards?
Big Data predicted that 67.53%
of you would remember!"



"You can't keep adjusting the data
to prove that you would be the best
Valentine's date for Scarlett Johansson."



"I don't like the look of this.
Searches for gravy and turkey stuffing
are going through the roof!"

DATA SCIENCE EVERYWHERE!...



"You can't keep adjusting the data to prove that you would be the best Valentine's date for Scarlett Johansson."

DATA SCIENCE EVERYWHERE!...

DO WE HAVE ANY ACTIONABLE ANALYTICS FROM OUR BIG DATA IN THE CLOUD?

oonist@gmail.com

YES, THE DATA SHOWS THAT MY PRODUCTIVITY PLUNGES WHENEVER YOU LEARN NEW JARGON.

nc./Dist. by Universal Uclick

MAYBE IN-MEMORY COMPUTING WILL ACCELERATE YOUR APPLICATIONS.

KDnuggets cartoon
© 2013 Ted Goff

WHAT DOES THE DATA TELL US TO DO?

WE ONLY HAVE BAD DATA ON THIS.

DILBERT.COM @SCOTTADAMSSAYS

DOES THE BAD DATA SUGGEST WE SHOULD DO WHAT WE WANTED TO DO ANYWAY?

WELL, YES.

4-3-18 ©2018 Scott Adams, Inc./Dist. by Andrews McMeel

THAT'S CALLED "GOOD DATA."

Valentine's date for Scarlett Johansson."



WHAT IS DATA SCIENCE?

WHAT IS DATA SCIENCE?

- “**Data science**, also known as **data-driven science**, is an interdisciplinary field of scientific methods, processes, algorithms and systems to extract knowledge or insights from data in various forms, either structured or unstructured, similar to data mining.”



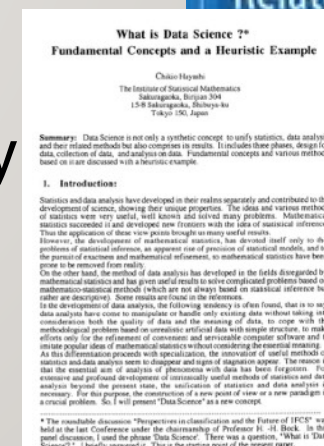
WIKIPEDIA
The Free Encyclopedia

WHAT IS DATA SCIENCE?

- “**Data science**, also known as **data-driven science**, is an interdisciplinary field of scientific methods, processes, algorithms and systems to extract knowledge or insights from data in various forms, either structured or unstructured, similar to data mining.”
- “Data science intends to analyze and understand actual phenomena with ‘data’. In other words, the aim of data science is to reveal the features or the hidden structure of complicated natural, human, and social phenomena with data from a different point of view from the established or traditional theory and method.”

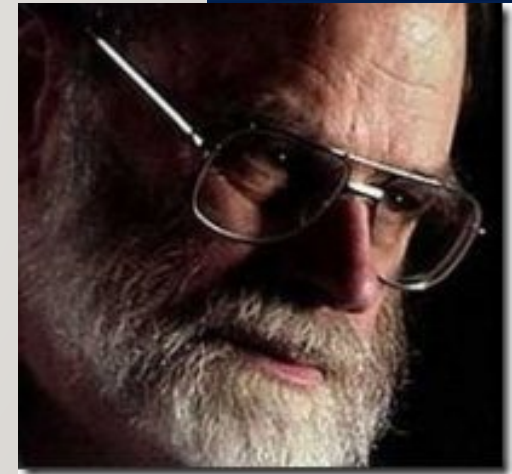
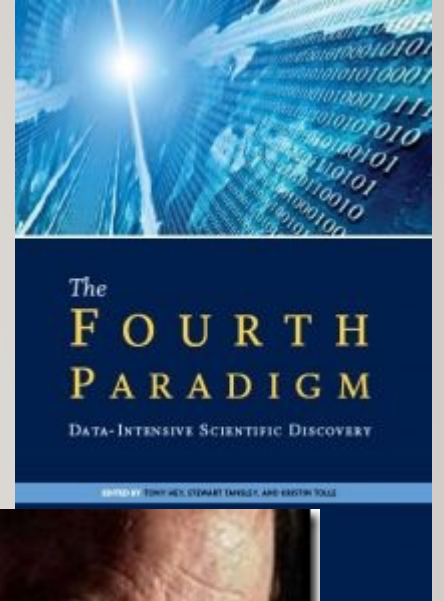


WIKIPEDIA
The Free Encyclopedia



WHAT IS DATA SCIENCE?

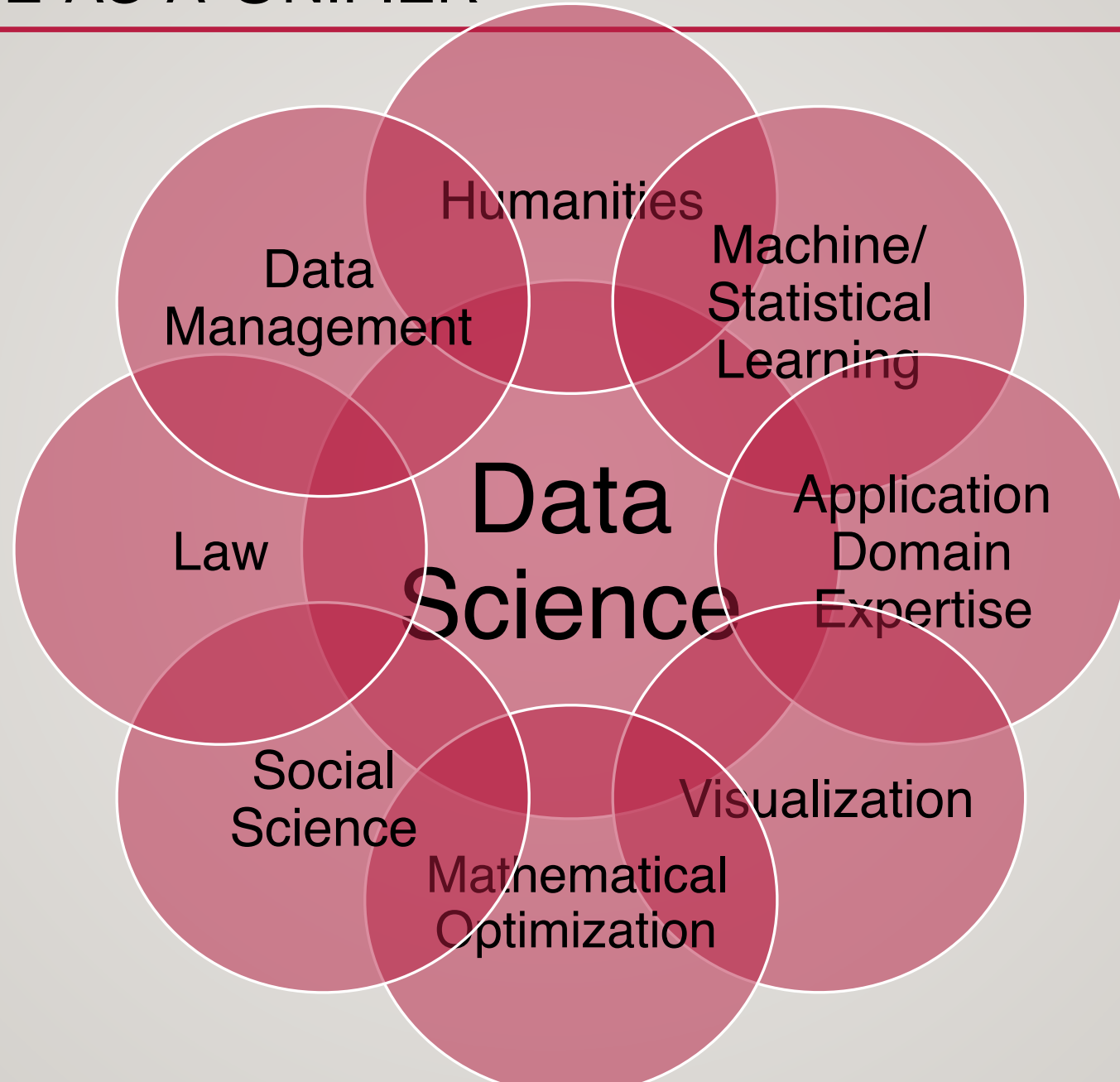
- Fourth paradigm
 - “... change of all sciences moving from observational, to theoretical, to computational and now to the 4th Paradigm – Data-Intensive Scientific Discovery”



WHAT IS IMPORTANT?

Need to solve a real problem using data...
No applications, no data science.

DATA SCIENCE AS A UNIFIER

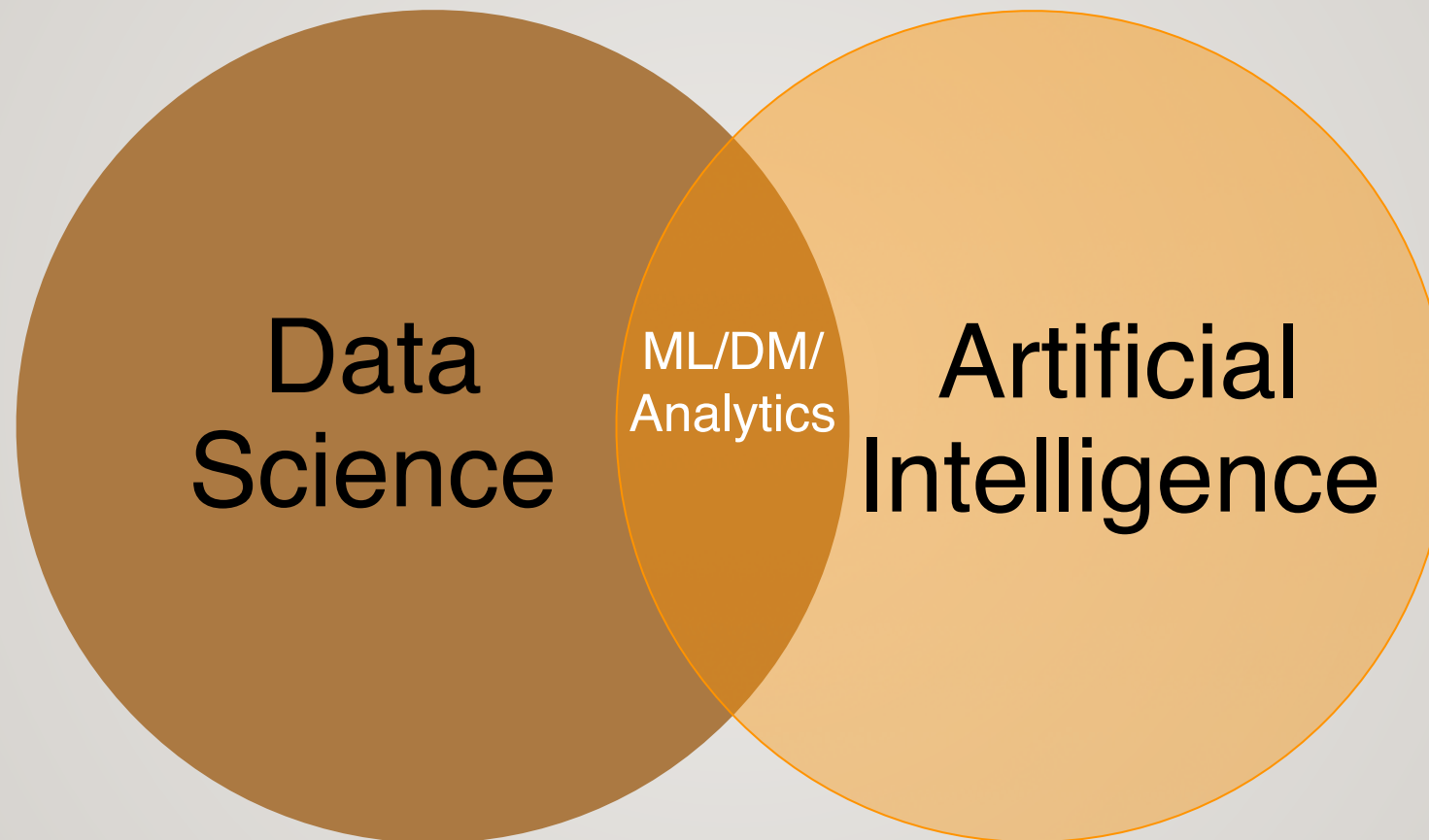


DATA SCIENCE AND BIG DATA

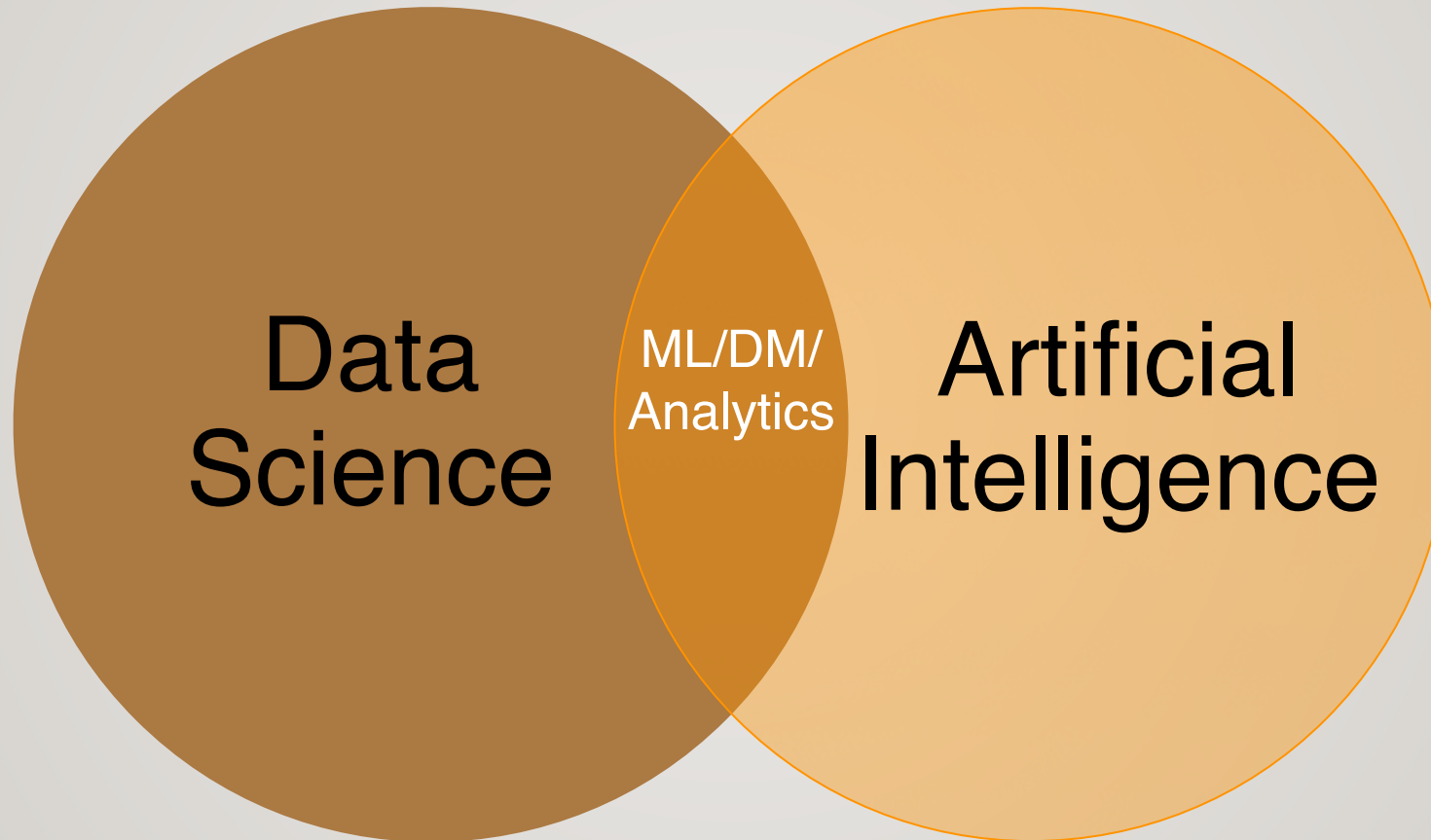
- They are not the “same thing”
- Big data = crude oil
 - Big data is about extracting “crude oil”, transporting it in “mega tankers”, siphoning it through “pipelines”, and storing it in “massive silos”
- Data science is about refining the “crude oil”

Carlos Samohano
Founder, Data Science London

DATA SCIENCE AND ARTIFICIAL INTELLIGENCE



DATA SCIENCE AND ARTIFICIAL INTELLIGENCE



**“Data science produces insights.
Machine learning produces predictions”**

DATA SCIENCE APPLICATION EXAMPLES

- Fraud detection
 - Investigate fraud patterns in past data
 - Early detection is important
 - Before damage propagates
 - Harder than late detection
 - Precision is important
 - False positive and false negative are both bad
 - Real-time analytics



DATA SCIENCE APPLICATION EXAMPLES

- Recommender systems
 - The ability to offer unique personalized service
 - Increase sales, click-through rates, conversions, ...
 - Netflix recommender system valued at \$1B per year
 - Amazon recommender system drives a 20-35% lift in sales annually
 - Collaborative filtering at scale



DATA SCIENCE APPLICATION EXAMPLES

- Predicting why patients are being readmitted
 - Reduce costs
 - Improve population health
 - Find the “why” behind specific populations being readmitted
 - Data lakes of multiple data sources
 - Investigate ties between readmission and socioeconomic data points, patient history, genetics, ...



DATA SCIENCE APPLICATION EXAMPLES

- “Smart cities”
 - Not well-defined



DATA SCIENCE APPLICATION EXAMPLES

- “Smart cities”
 - Not well-defined



illisphotography / Getty / Emily Jan / The Atlantic

Stop Saying 'Smart Cities'

Digital stardust won't magically make future cities more affordable or resilient.

BRUCE STERLING | FEB 12, 2018 | TECHNOLOGY

DATA SCIENCE APPLICATION EXAMPLES

- “Smart cities”
 - Not well-defined
 - Generally refers to using data ICT to
 - Better plan communities
 - Better manage assets
 - Reduce costs
 - Deploy open data to better en with community



lillisphotography / Getty / Emily Jan / The Atlantic

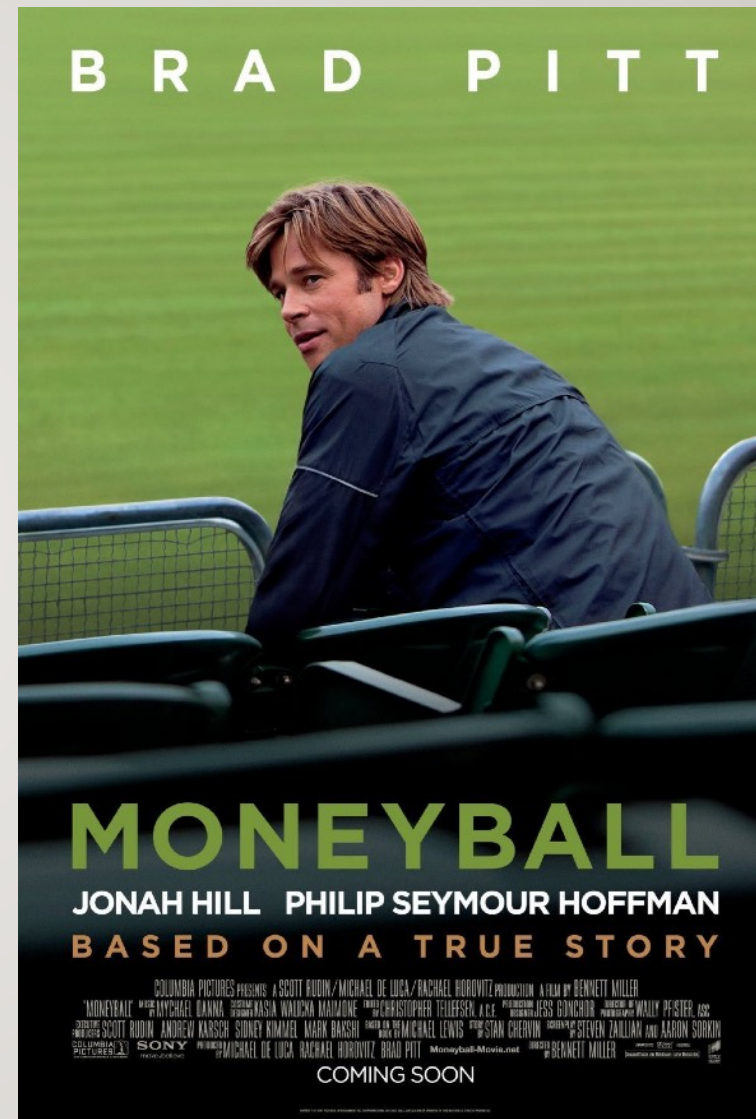
Stop Saying 'Smart Cities'

Digital stardust won't magically make future cities more affordable or resilient.

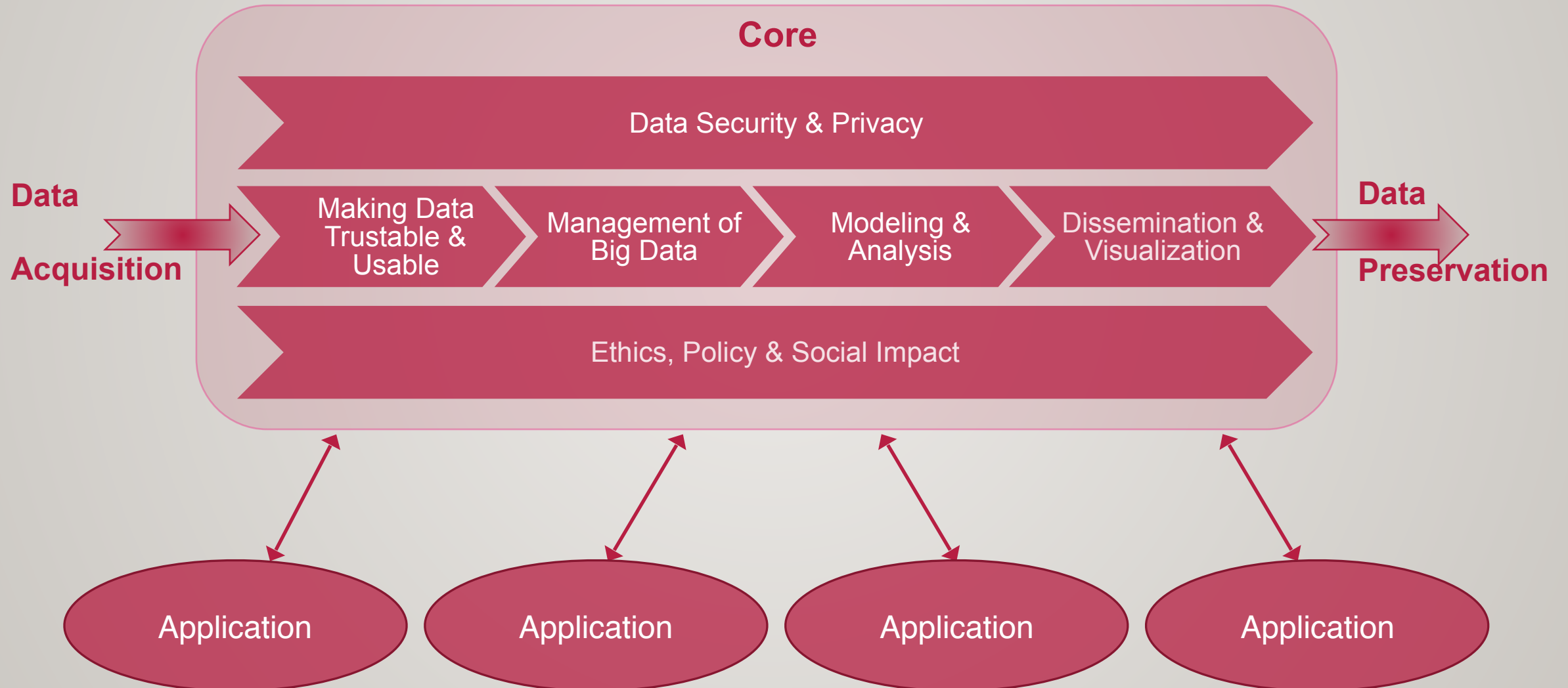
BRUCE STERLING | FEB 12, 2018 | TECHNOLOGY

DATA SCIENCE APPLICATION EXAMPLES

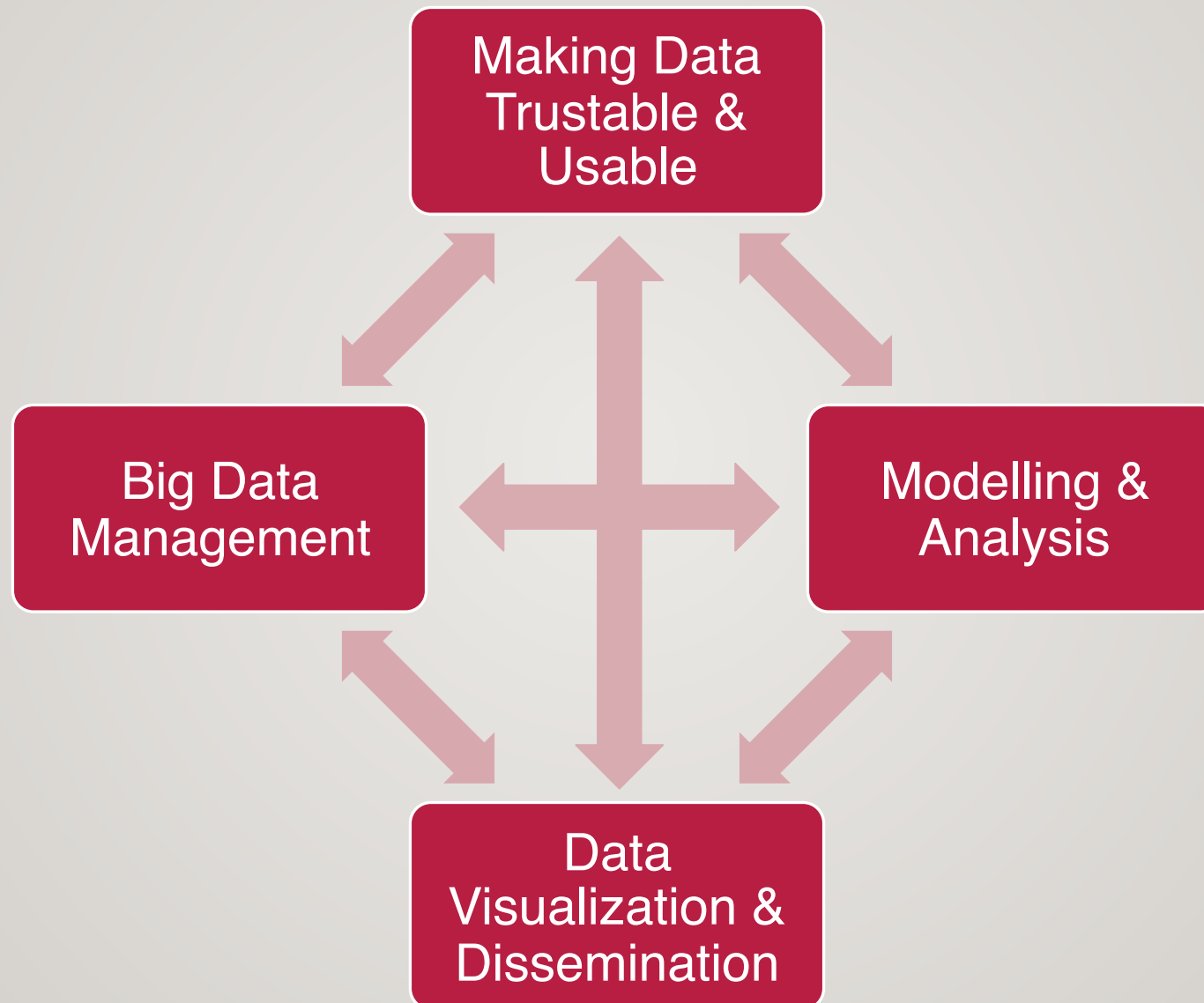
- Moneyball
 - How to build a baseball team on a very low budget by relying on data
 - *Sabermetrics*: the statistical analysis of baseball data to objectively evaluate performance
 - 2002 record of 103-59 was **joint best** in MLB
 - Team salary budget: \$40 million
 - Other team: Yankees
 - Team salary budget: \$120 million



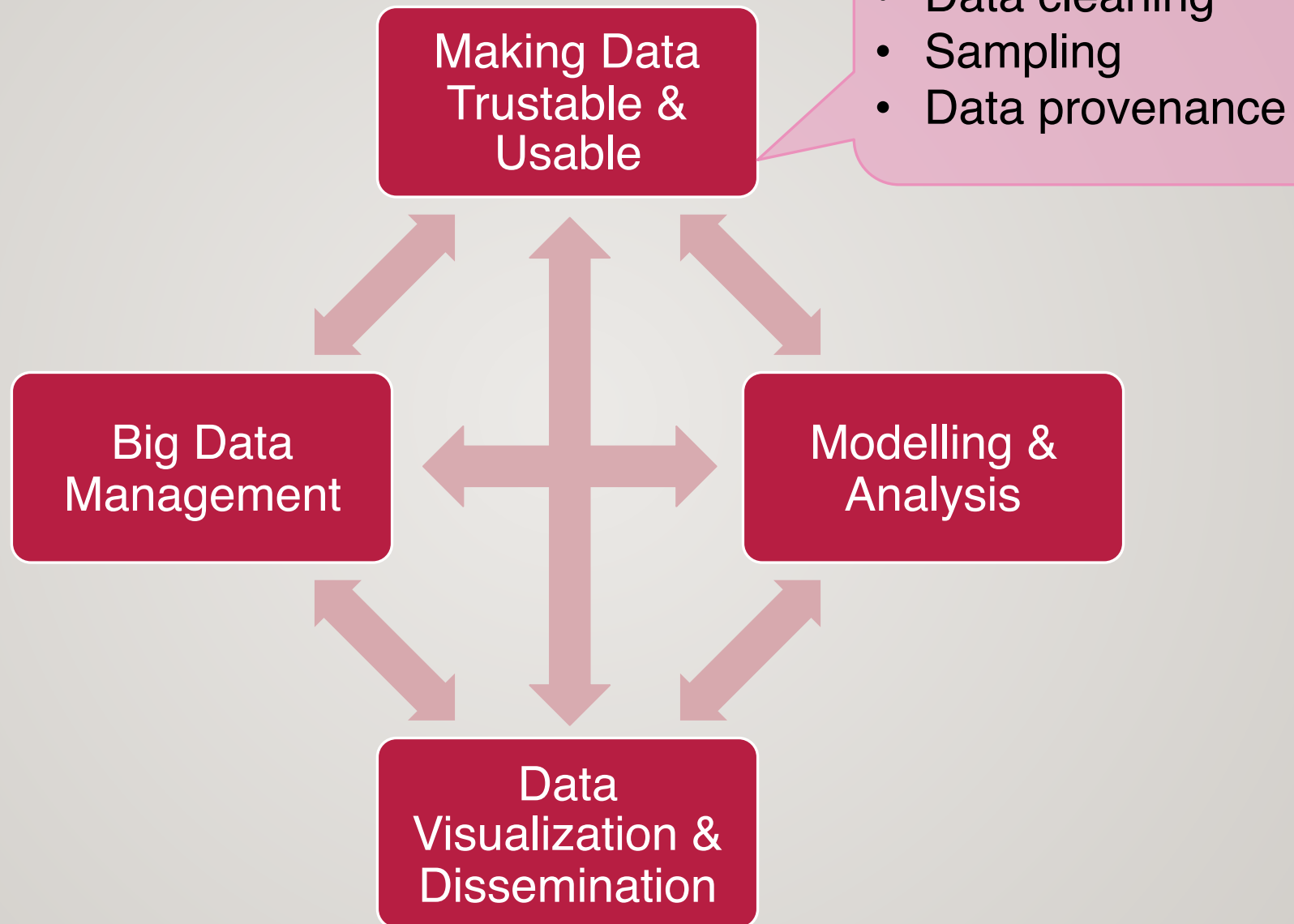
HOLISTIC APPROACH TO DATA SCIENCE



CORE RESEARCH ISSUES & INTERACTIONS



CORE RESEARCH ISSUES & INTERACTIONS



CORE RESEARCH ISSUES & INTERACTIONS

- Data lakes
- Batch & online access
- Platforms

Making Data
Trustable &
Usable

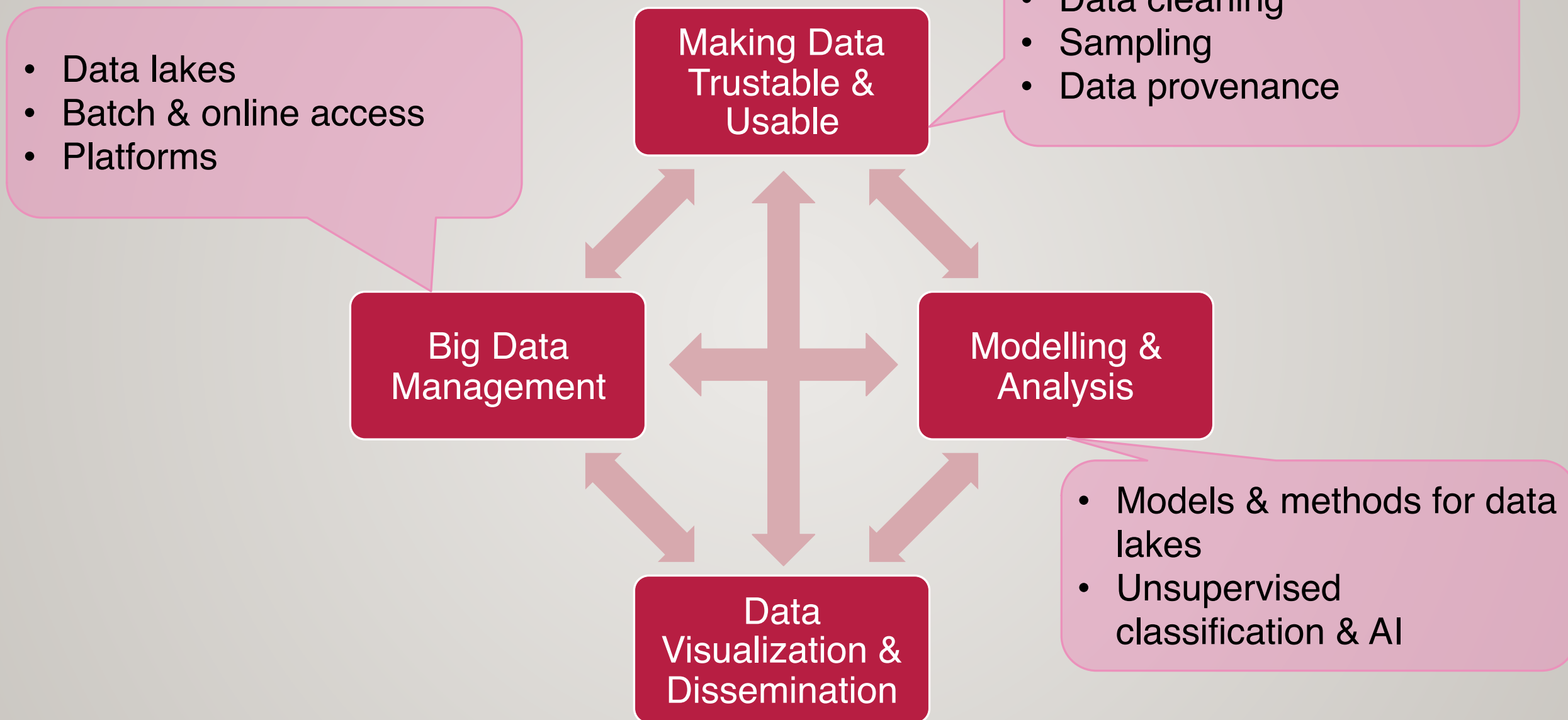
- Data cleaning
- Sampling
- Data provenance

Big Data
Management

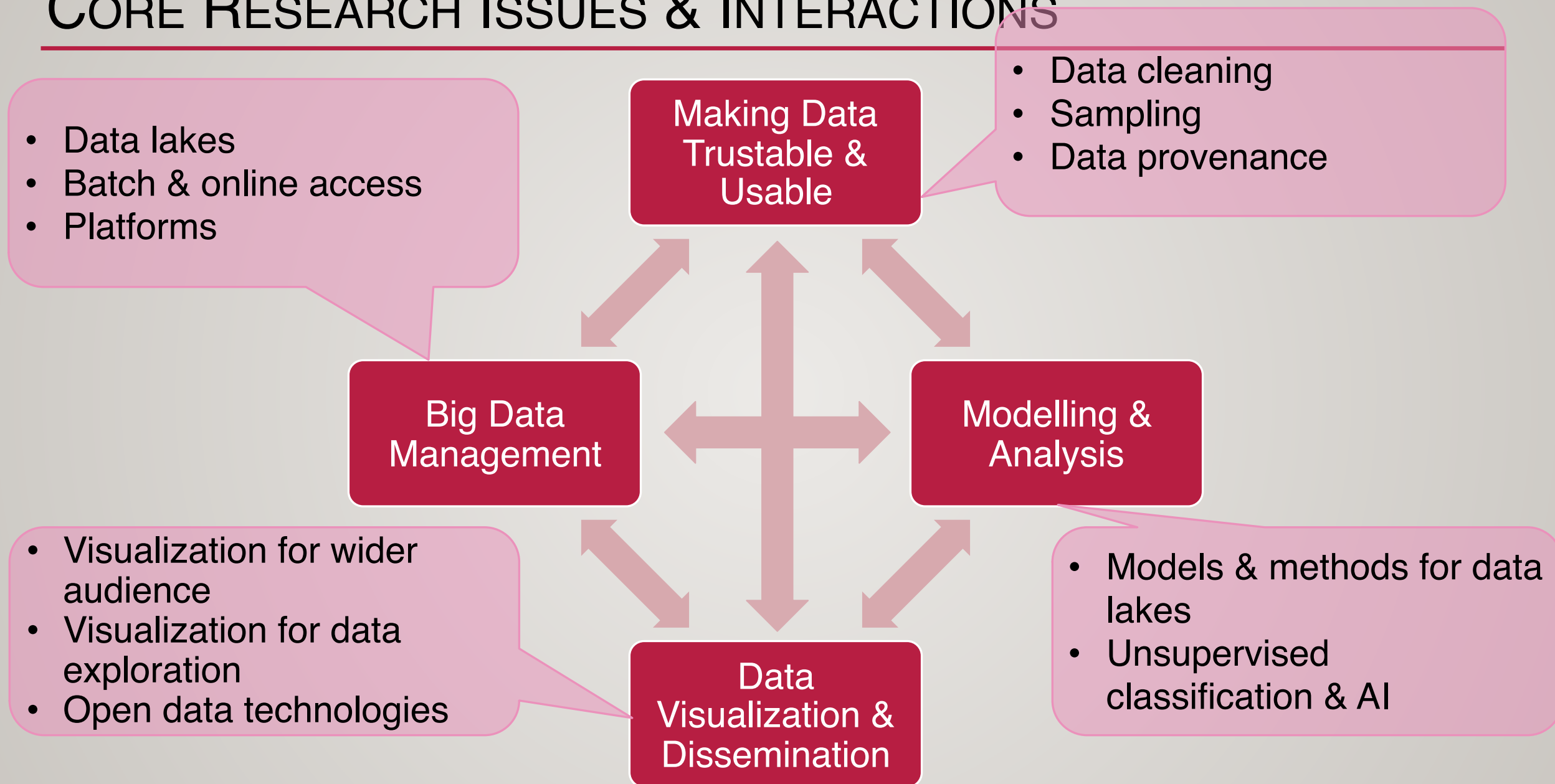
Modelling &
Analysis

Data
Visualization &
Dissemination

CORE RESEARCH ISSUES & INTERACTIONS



CORE RESEARCH ISSUES & INTERACTIONS



CORE RESEARCH ISSUES & INTERACTIONS

