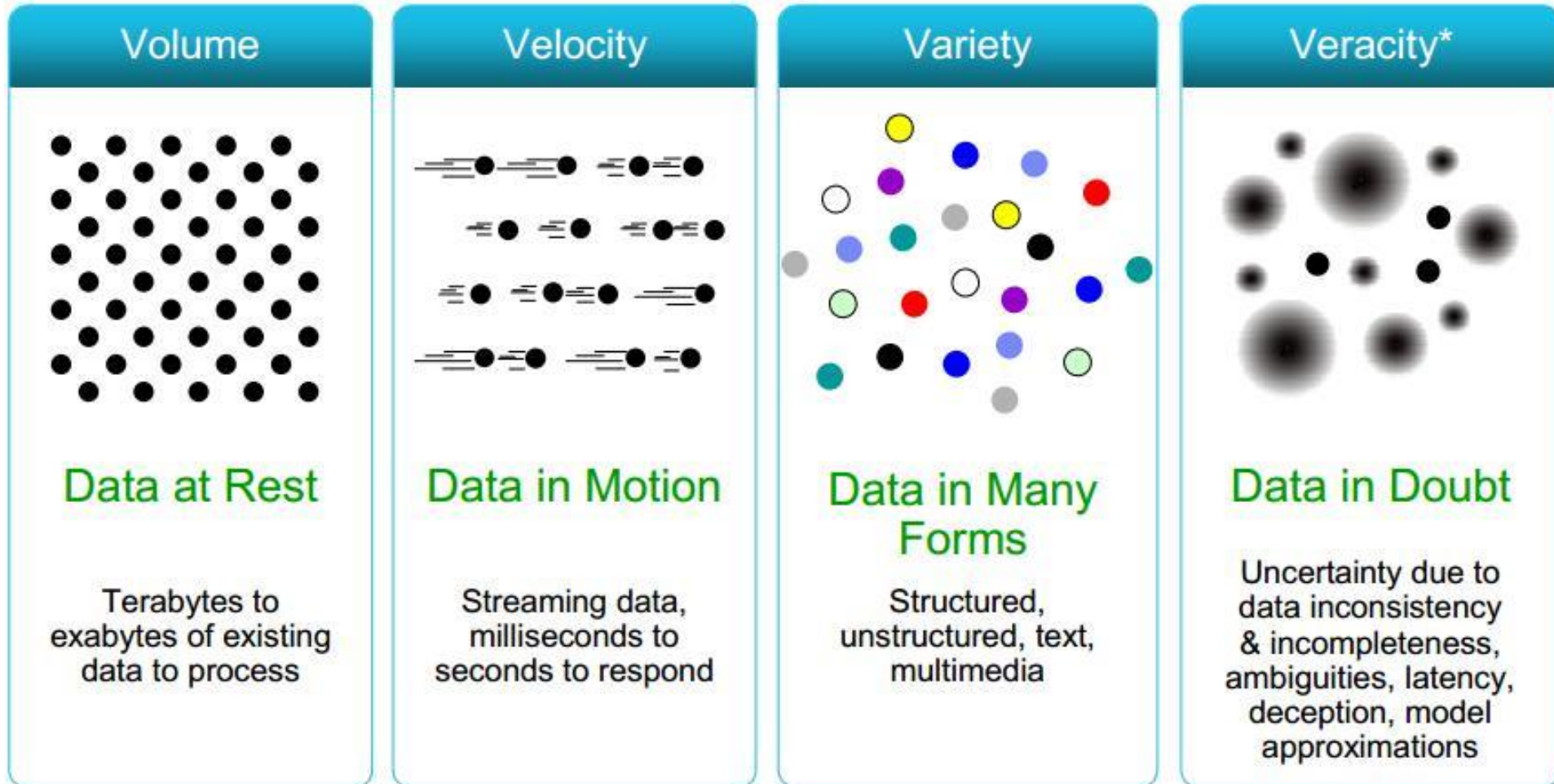


Module 1: Big Data

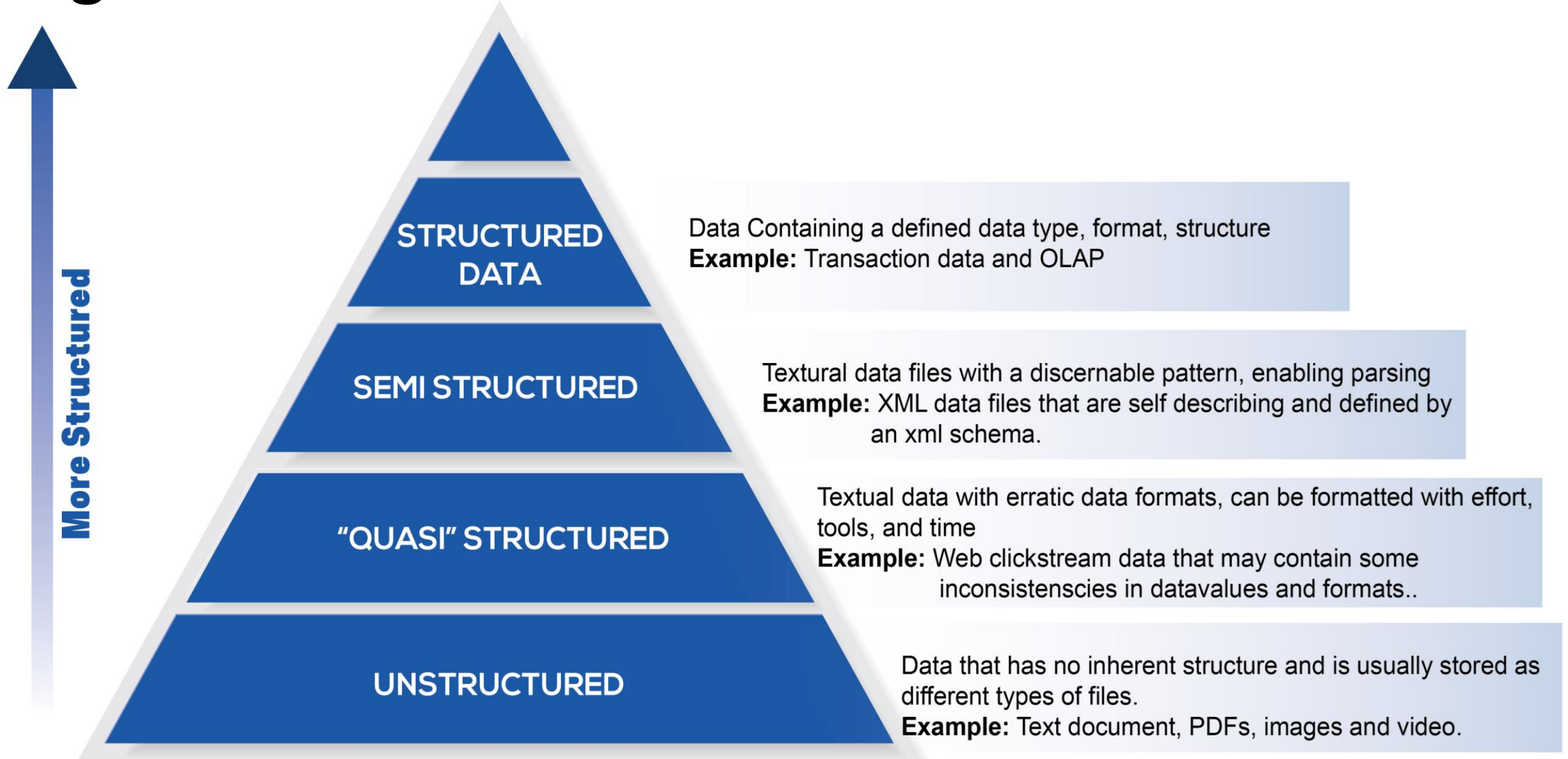


“Big data is a collection of data from traditional and digital sources inside and outside your company that represents a source for ongoing discovery and analysis.” (Lisa Arthur, CMO Network, 8/15/2013). “Big Data” is data whose scale, diversity, and complexity require new architecture, new tools, techniques, algorithms, and analytics to manage it and extract value and hidden knowledge from it...

Some Make it 4V's



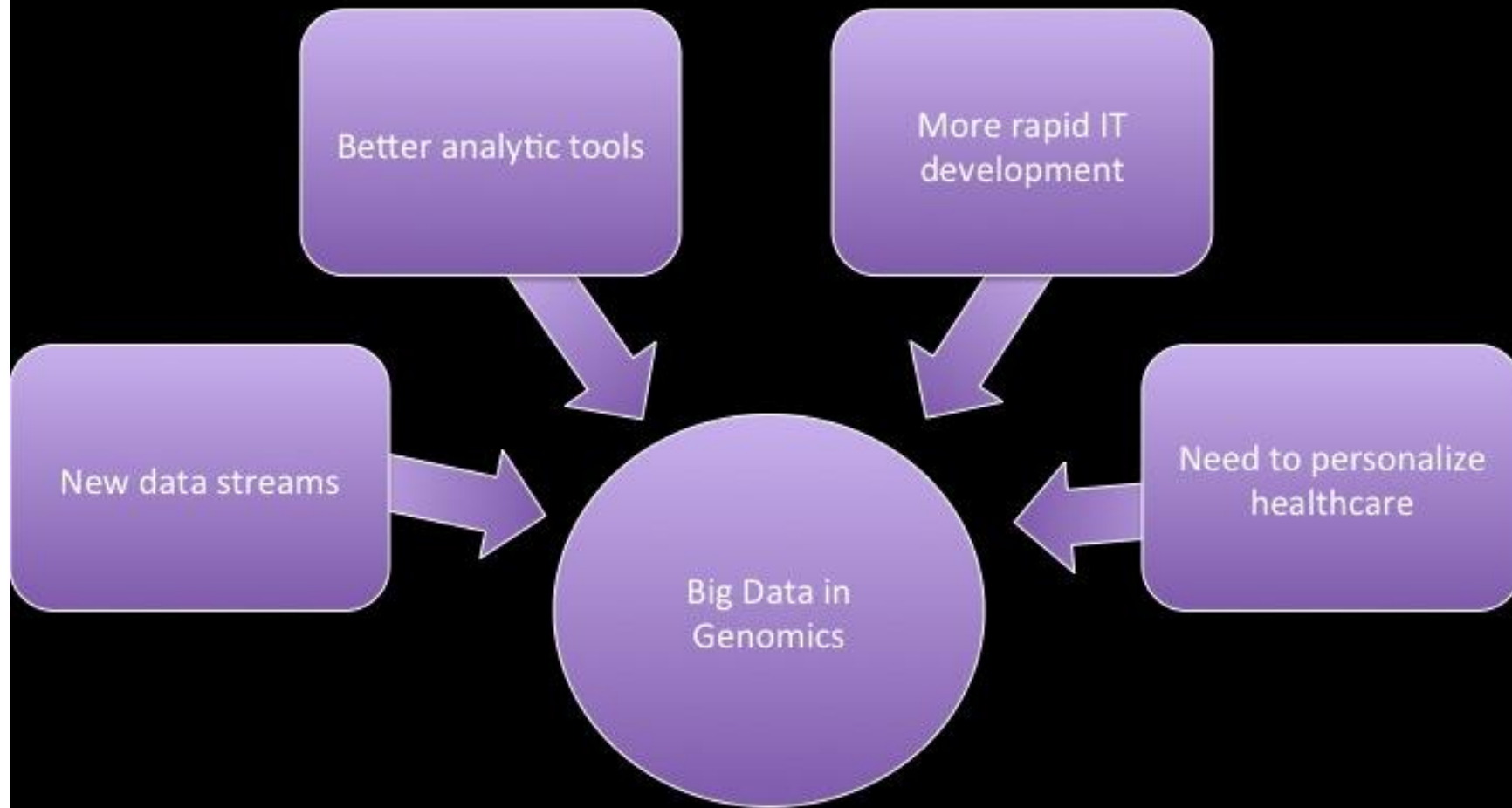
Big Data Characteristics: Data Structure

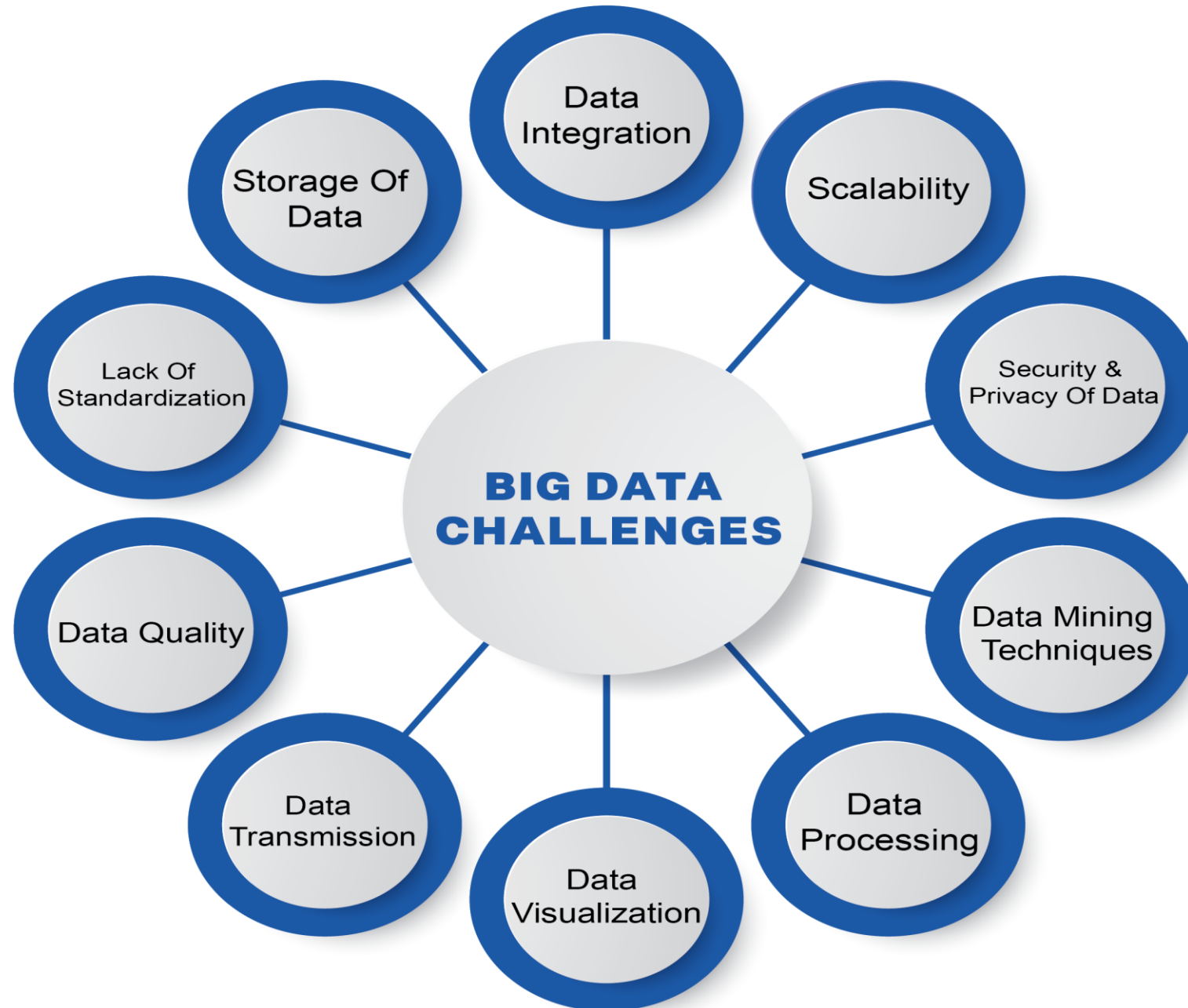


Traditional Data vs. Big Data

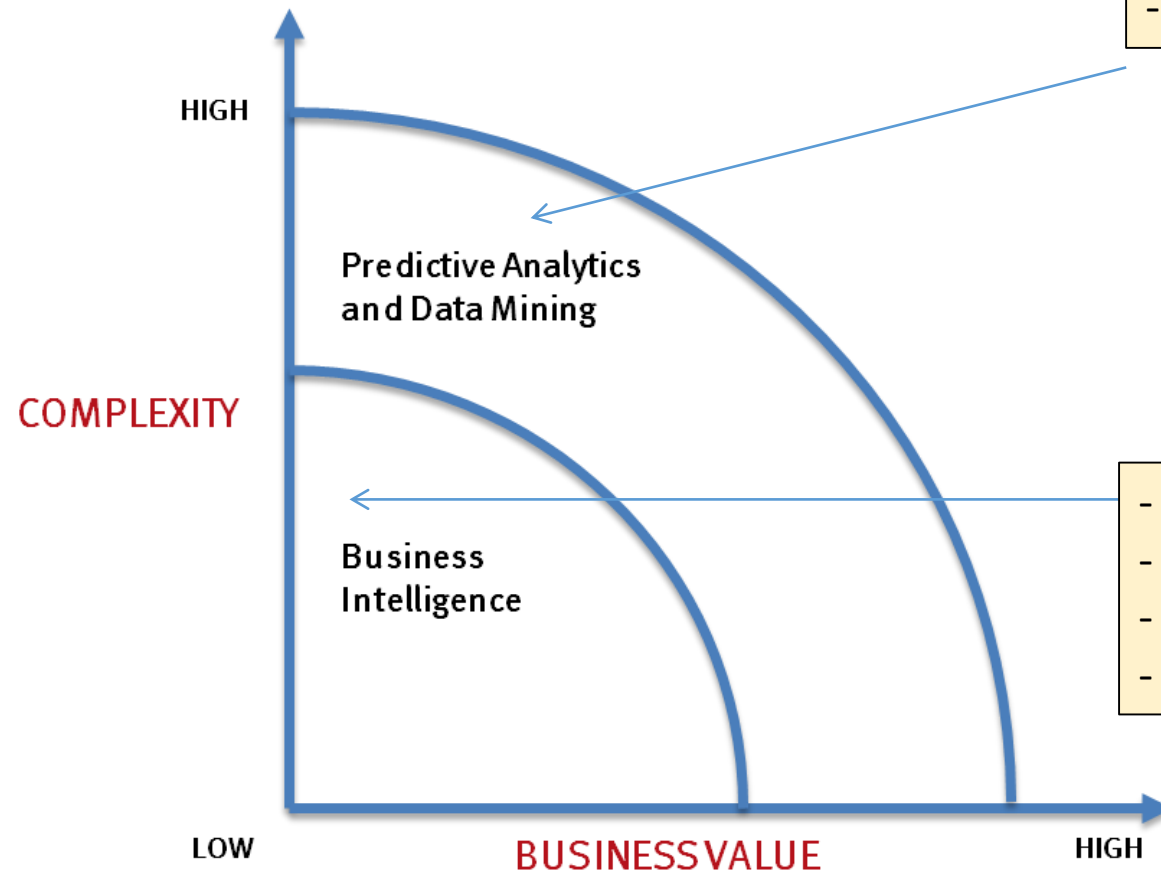
- Challenges
- Advantages

Why Big Data in Genomics Now?





What's Driving Big Data



- Optimizations and predictive analytics
- Complex statistical analysis
- All types of data, and many sources
- Very large datasets
- More of a real-time

Common questions

- What if...
- What's the optimal scenario for our business?
- What will happen next?
What if these trends continue? Why is this happening?

- Ad-hoc querying and reporting
- Data mining techniques
- Structured data, typical sources
- Small to mid-size datasets

Common questions

- What happened last quarter?
- How many did we sell?
- Where is the problem? In which situation?

DATA PREPARATION

DATA CLEANING

INCONSISTENT DATATYPES
MISSPELLED ATTRIBUTES
MISSING AND DUPLICATE VALUES

TRANSFORMATION



EXPLORATORY DATA ANALYSIS



DEFINES AND REFINES
THE SELECTION OF FEATURE
VARIABLES THAT WILL BE USED
IN THE MODEL DEVELOPMENT

DATA MODELING

simplilearn

KNN



NAIVE BAYES

DECISION TREE

VISUALIZATION AND COMMUNICATION



WHAT IS DATA SCIENCE?

DATA ACQUISITION

- WEB SERVERS
- LOGS
- DATABASES
- APIS
- ONLINE REPOSITORIES



WHY?...WHY?...WHY?....



DEPLOYS AND

