#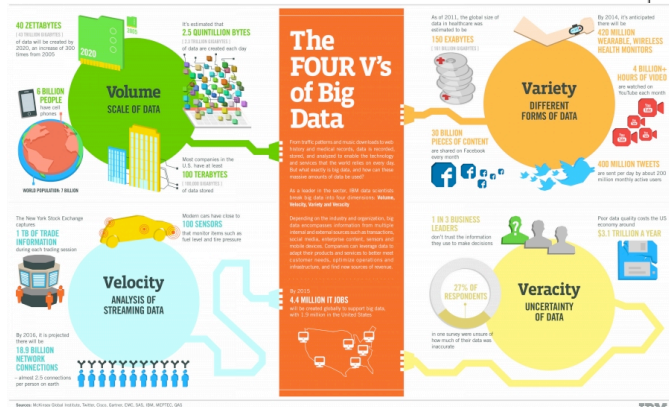 Linking collections to related resources: Multi-scale, multi-dimensional, multi-disciplinary collaborative research in biodiversity.  Is this a "Big Data" paradigm?

*Reed Beaman,*
*, University of Florida, Gainesville, FL, USA*

# The 3 or 4 Vs of Big Data



The FOUR V's of Big Data

**Volume** — SCALE OF DATA
**Velocity** — ANALYSIS OF STREAMING DATA
**Variety** — DIFFERENT FORMS OF DATA
**Veracity** — UNCERTAINTY OF DATA

Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS
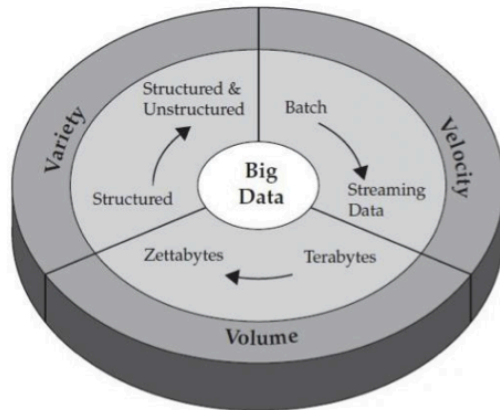


Figure 1-1  IBM characterizes Big Data by its volume, velocity, and variety—or simply, V³.



Figure 1 — Data Management Solutions

**Volume**
 - Tiered storage/hub and spoke
 - Selective data retention
 - Statistical sampling
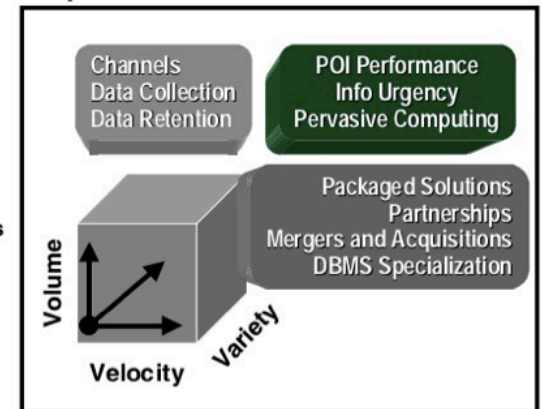 - Redundancy elimination
 - Offload "cold" data
 - Outsourcing

**Velocity**
 - Operational data stores
 - Data caches
 - Point-to-point data routing
 - Balance data latency with decision cycles

**Variety**
 - Inconsistency resolution
 - XML-based "universal" translation
 - Application-aware EAI adapters
 - Data access middleware and ETLM
 - Distributed query management
 - Metadata management

**E-Business-Driven Information Explosion Factors**

Channels / Data Collection / Data Retention

POI Performance / Info Urgency / Pervasive Computing

Packaged Solutions / Partnerships / Mergers and Acquisitions / DBMS Specialization
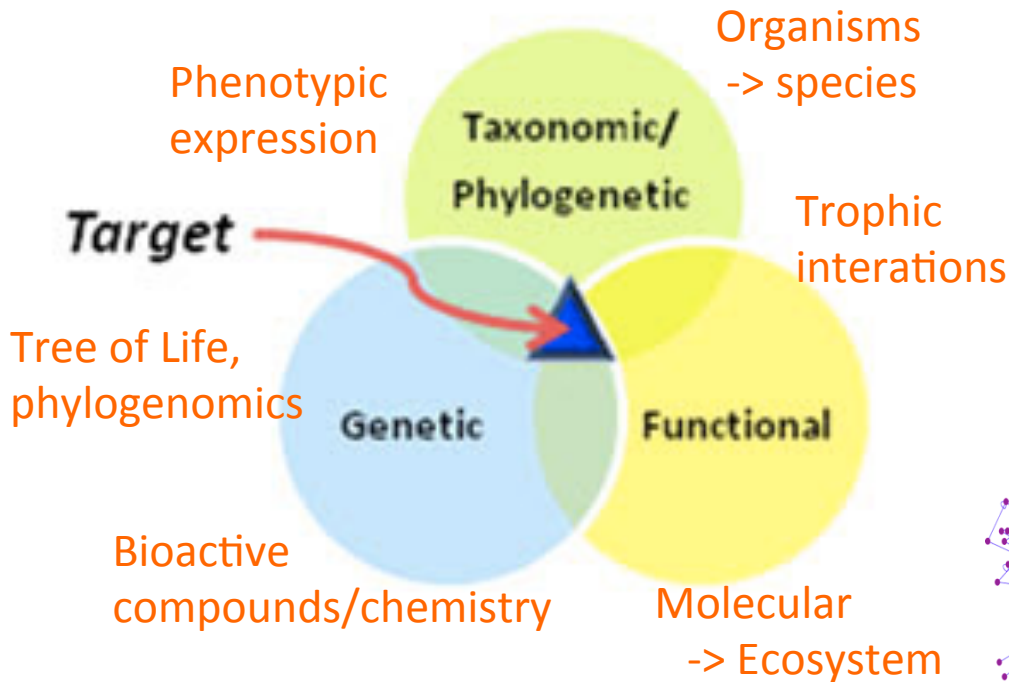
Volume / Velocity / Variety

*Extending data management options enables greater returns on information assets*
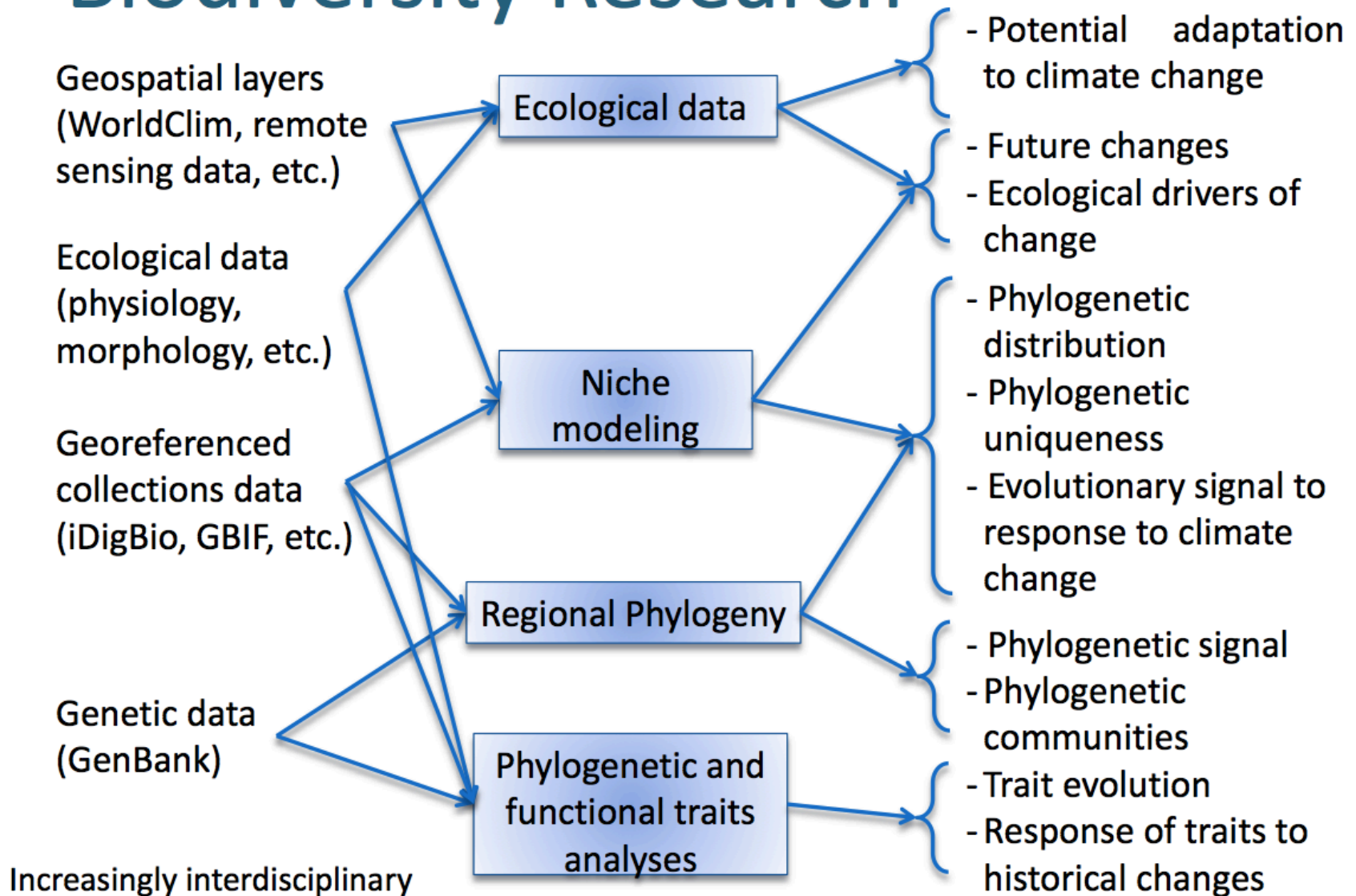
Source: META Group

"Big data is data that's an order of magnitude bigger than you're accustomed to, Grasshopper."   Doug Laney, Gartner

# Integrative Biodiversity: Multiscale, Multi-disciplinary

- US NSF Dimensions of Biodiversity program)
  - Interaction at the intersection of taxonomic, genetic, functional domains



Phenotypic expression

Organisms -> species

Taxonomic/ Phylogenetic

Target

Trophic interations

Tree of Life, phylogenomics

Genetic

Functional

Bioactive compounds/chemistry

Molecular -> Ecosystem

# Biodiversity Research

Geospatial layers (WorldClim, remote sensing data, etc.)

Ecological data (physiology, morphology, etc.)

Georeferenced collections data (iDigBio, GBIF, etc.)

Genetic data (GenBank)

Increasingly interdisciplinary

**Ecological data**

**Niche modeling**

**Regional Phylogeny**

**Phylogenetic and functional traits analyses**

- Potential    adaptation to climate change

- Future changes
- Ecological drivers of change

- Phylogenetic distribution
- Phylogenetic uniqueness
- Evolutionary signal to response to climate change

- Phylogenetic signal
- Phylogenetic communities

- Trait evolution
- Response of traits to historical changes

Source:  Andrea Matsunaga

Big data is a given for genomics, high throughput sequencing, analysis, and visualization

What about all the other data that **relates** to genetic and genomic data?

# 4Vs for Biodiversity Big data

- Volume: billion or more specimens, 2-10 million species (excluding microbial), 10 billion plus related edges
- Velocity: Snail's pace? 250 year long-tail legacy of taxonomic data -> rapid digitization <-> large scale genomic sequencing
- Variety: Occurrences, sequences, morphological, geospatial; structured and unstructured
- Veracity: Very challenging to validate?

# Figure 3. Linking samples and derivatives from the Moorea Biocode project.

BiSciCol (Biological Science Collections Tracker) use case:

Every specimen links to a multitude of parent and derivative data. Users of biodiversity data need to be able to *easily and quickly* see these relationships

# The "Big" in Ecological Big Data

*The defining aspect of ecological Big Data is not raw size but another dimension: complexity.*

Dave Schimel,
(former) NEON
Chief Scientist



Modified from the image of Nicolle Rager Fuller, National Science Foundation, 2007

# 4Cs of Biodiversity Big Data

- Complexity:  scale, interactions (e.g., food webs) between individuals, populations, species, environments (cf. story lines)

- Collaboration:  International and multidisciplinary

- Citizen Science:  Increasing as a solution to digitization

- Completeness:  Will we always be 10% complete, and can we validate and create the linkages?

# Figure 2. Core terms of the Biological Collections Ontology (BCO) and their relations to upper ontologies.

# Software Defined System

Adjusts to changing needs and environments

From Virtual Machines to Virtual Clusters

Application VM

Virtual Cluster

Overlay Network

Data sharing over multiple networks

Data Server

Other Networks

Move the software to the data

# Trust Envelope

AIST Japan may have more compute resources

LifeMapper AIST

LifeMapper Virtual Cluster

Satellite imagery

Sensitive biodiversity data and UAV (Drone) imagery

Overlay Network

Sensitive or licensed data may not be portable

iDigBio, GBIF

# Integrative Biodiversity

- Grand challenge science:  Big data is about asking big questions