Paralisin: Tghty copied Loosely copied Applications Implications Implications Implications Operating Windows Heterogeneous Implications Implications Virtualization Implications Implications Implications Implications Implications Virtualization Implications Implications Implications Implications Implications Implications Virtualization Implications Implications Implications Implications Implications Implications Data center mechanical and electrical Implications Implications Implications Implications Nature Reviews Implications Implications Implications Implications Implications Networking Implications Implications Implications Implications Implications Networking Implications Implications Implications Implications Implications Matching Implications Implications Implications Implications Implications Matching Implications Implications

Cloud computing

BINF 630: Bioinformatics Methods

Iosif Vaisman

Email: ivaisman@gmu.edu

Schadt et al., 2010

Cloud computing and DNA sequencing



Types of computational environments





Types of computational environments

Environment	URL
Cloud computing	
Amazon Elastic Compute Cloud	http://aws.amazon.com/ec2
Bionimbus	http://www.bionimbus.org
NSF CluE	http://www.nsf.gov/cise/clue/index.jsp
Rackspace	http://www.rackspacecloud.com
Science Clouds	http://www.scienceclouds.org
Heterogeneous computing	1
NVIDIA GPUs	http://www.nvidia.com
AMD/ATI GPUs	http://www.amd.com
Heterogeneous cloud com	puting
SGI Cyclone Cloud	http://www.sgi.com/products/hpc_cloud/cyclone
Penguin Computing On Demand	http://www.penguincomputing.com/POD/Summary
GPU araphics processing unit	NSE US National Science Foundation

Schadt et al., 2010

Infrastructure organization

Defining Big Data

NOT JUST SIZE

(Mbs

The three Vs of Big Data: volume, variety and velocity (D.Lanev, 2001)

Elements of "Big Data" include:

•The degree of complexity within the data set

•The amount of value that can be derived from innovative vs. non-innovative analysis techniques

•The use of longitudinal information supplements the analysis

http://mike2.openmethodology.org/wiki/Big_Data_Definition

L D Stein, 2010

Genome sequencing costs





Sequencing and storage cost

Knowledge Discovery

Knowledge is a.pattern that exceeds certain threshold of interestingness.

Factors that contribute to interestingness:

coverage confidence statistical significance simplicity unexpectedness actionability

Data Mining

- Data mining is the exploration and analysis, by automatic or semiautomatic means, of large quantities of data in order to discover meaningful patterns and rules
- Common data mining tasks
 - Classification
 - Estimation
 - Prediction
 - Affinity Grouping
 - Clustering
 - Description

Knowledge Discovery

- · Directed and Undirected KD
- · Directed KD
 - Purpose: Explain value of some field in terms of all the others
 - Method: We select the target field based on some hypothesis about the data. We ask the algorithm to tell us how to predict or classify it
 - Similar to hypothesis testing (e.g., in regression modeling) in statistics

Knowledge Discovery

- Undirected KD
 - Purpose: Find patterns in the data that may be interesting
 - Method: clustering, affinity grouping
 - Closest to ideas of machine learning in artificial intelligence
- · Comparison
 - UKD helps us to recognize relationships & DKD helps us to explain them

Classification

Classifying observations into different categories given characteristics

Estimation

• Rules that explain how to estimate a value given characteristics

Prediction

 Rules that explain how to predict a future value or classification, given characteristics

Affinity Grouping

Grouping by relations (not by characteristics)

Clustering

- Segmenting a diverse population into more similar groups
- In clustering, there are no pre-defined classes and no examples. Records are grouped together by some similarity measure.

Knowledge Discovery

