



A Platform-as-a-Service for Biobanking

www.biobankcloud.com

- **Presented by: Paulo Esteves Veríssimo**
- **Univ. de Lisboa, Faculdade de Ciências (FCUL), LaSIGE, Portugal,**
- pjv@di.fc.ul.pt <http://www.di.fc.ul.pt/~pjv>
- **64th IFIP WG 10.4 Meeting, Visegrad - Hungary**
- **June 27-30, 2013**



Financed by the European Commission 7th Framework Programme.



BiobankCloud – a PaaS for Biobanking

Duration: Dez 2012-2015

Budget: ~3Mio Euros



A Biobank is :
"biological material from one or several human beings collected and stored indefinitely or for a specified time and whose origin can be traced to the human or humans from whom it originates"

Goal

Provide a viable platform-as-a-service (PAAS) for Biobanking to support the secure storage, analysis, and sharing of genomic data.

Research Challenges

Design and build services for cloud platforms that securely store and analyze sensitive personal data, while conforming to a regulatory framework.

Genomics and Big Data

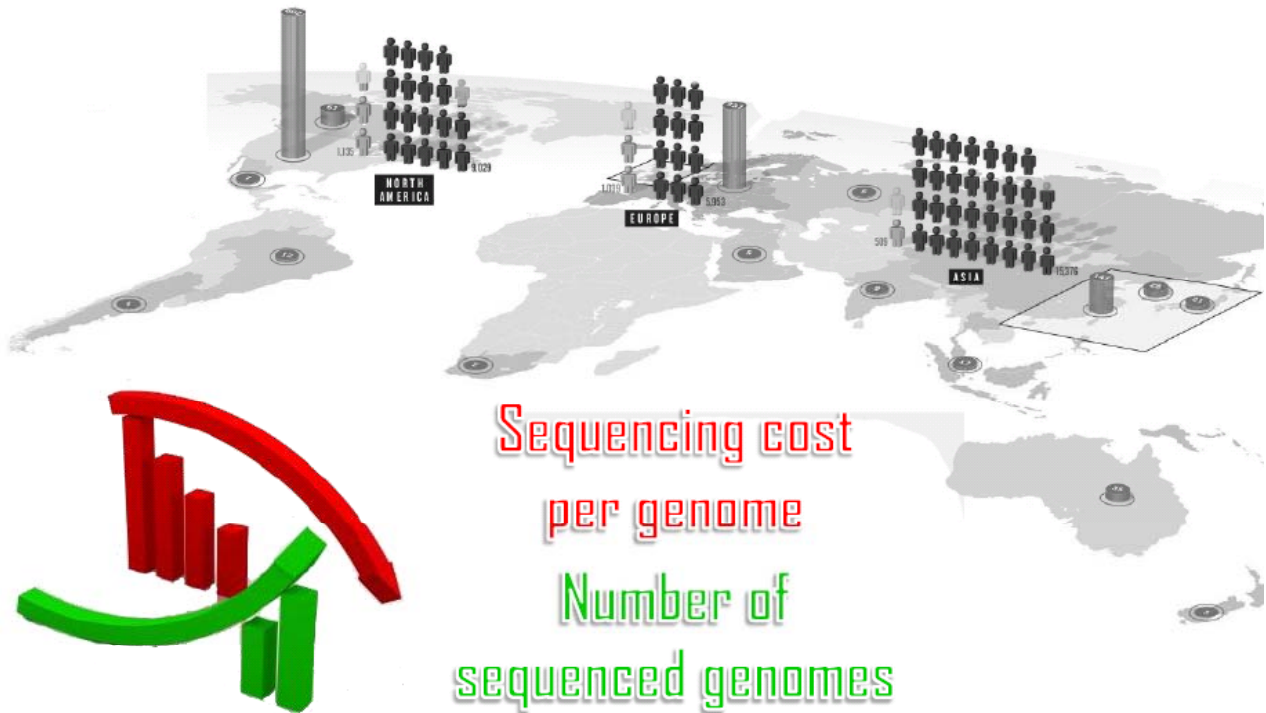
- We will soon be generating more genomic data than we can currently securely store and process
- Handling terabytes of information has become the norm for genomics, but it's still a challenge to set up.
- The research community needs systems that are efficient and secure.

Urgent need for better systems to store and process genomic data

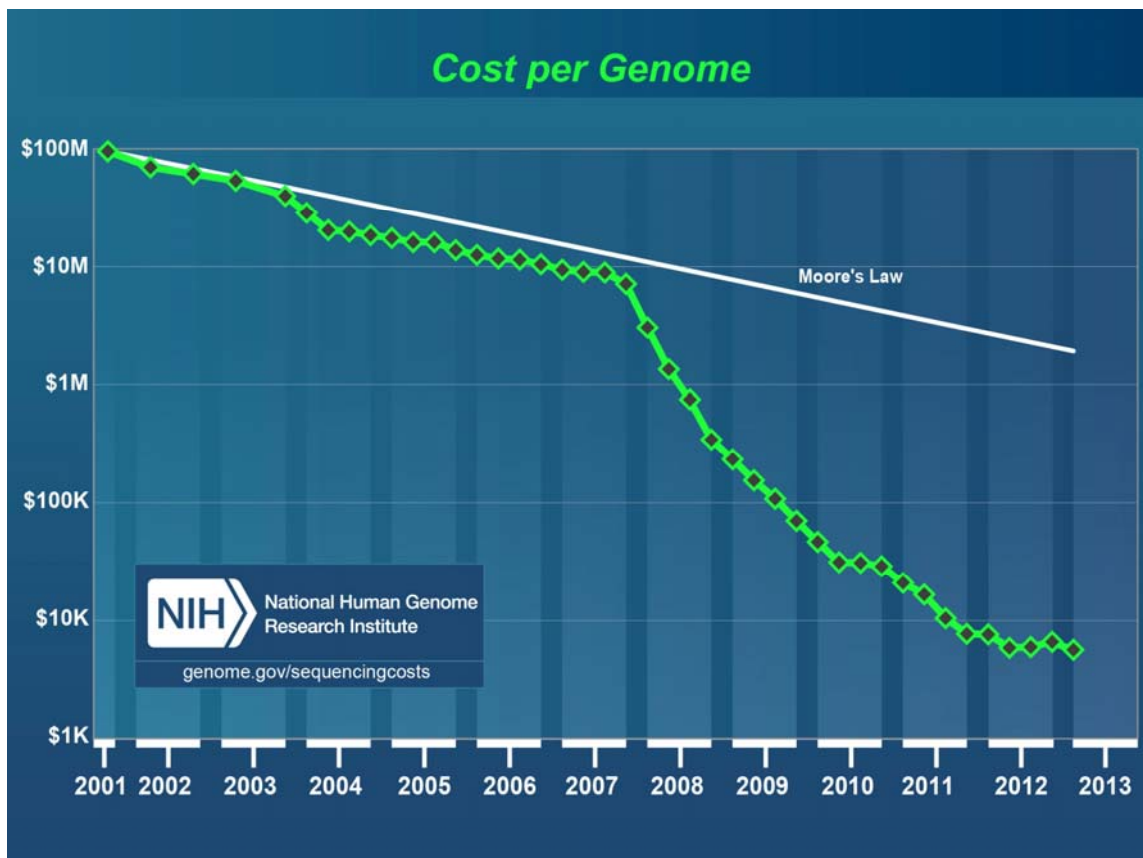


A scent of the future trends ...

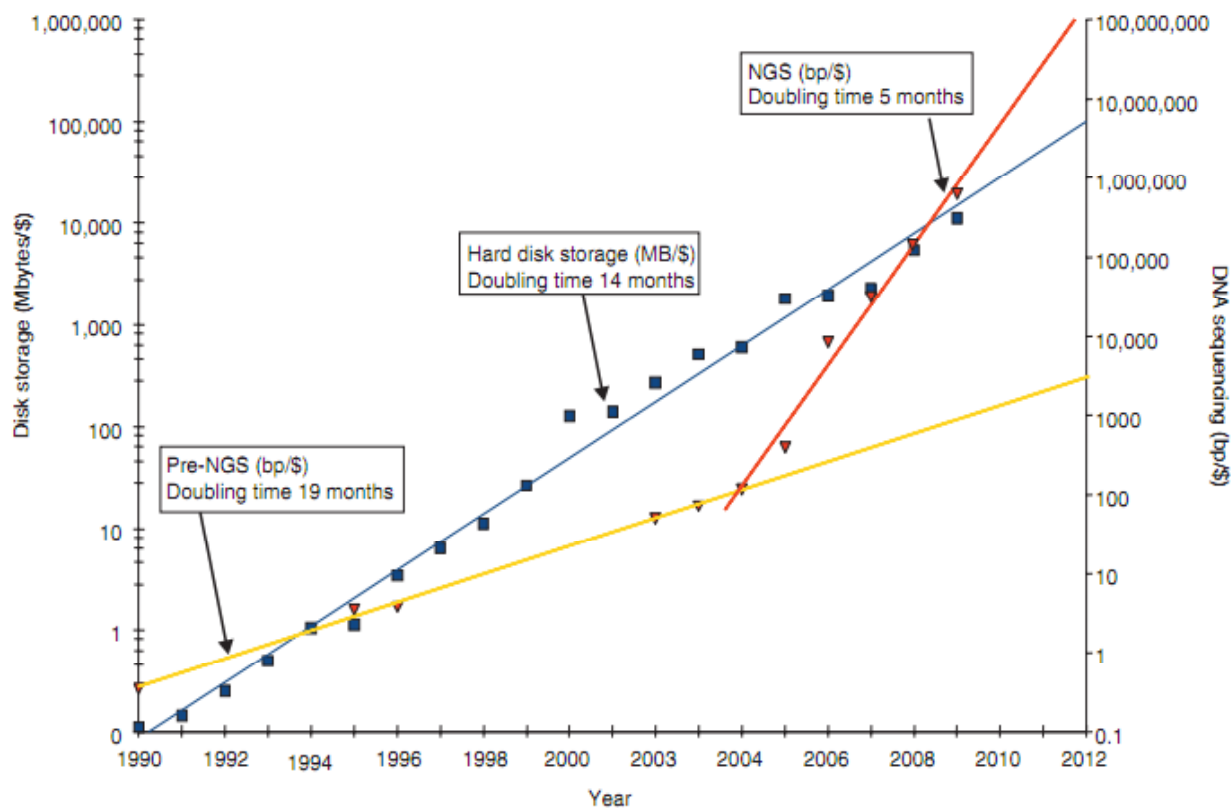
The cost per genome is bound to decrease



Cost per genome vs. Moore's Law



Growth in NGS capability vs. Kryder's law



- Growth in disk storage capacity vs. growth in NGS data generation capability

Current reality already significant ...

Increasing number of research databases

1000 Genomes

A Deep Catalog of Human Genetic Variation



Complete
genomics
A BGI Company



GenBank



UK
10K

RARE GENETIC VARIANTS IN HEALTH AND DISEASE



GENOME 10K



Personal
Genome
Project

LifeGene: ~ 500,000 people



Sweden has 73 **Opt-Out** “Quality Registries”

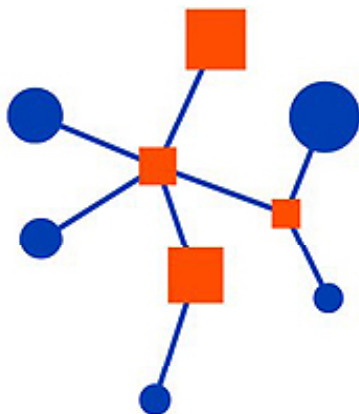


Nationella Kvalitetsregisteret

http://www.kvalitetsregister.se/om_kvalitetsregister/quality_registries

A Coalition of the Willing

- In Europe, BBMRI is a trans-national network of Biobanks with **the goal of sharing study data** for research



BBMRI

Biobanking and
Biomolecular
Resources Research
Infrastructure

Storing and computing with this data ...

Scalable storage on commodity hardware

- Many genomics research groups still run expensive SANs (storage area networks) or file systems that require expensive interconnects.



Reliable Storage Service



Unreliable Commodity Servers

Immutable Genomic Data

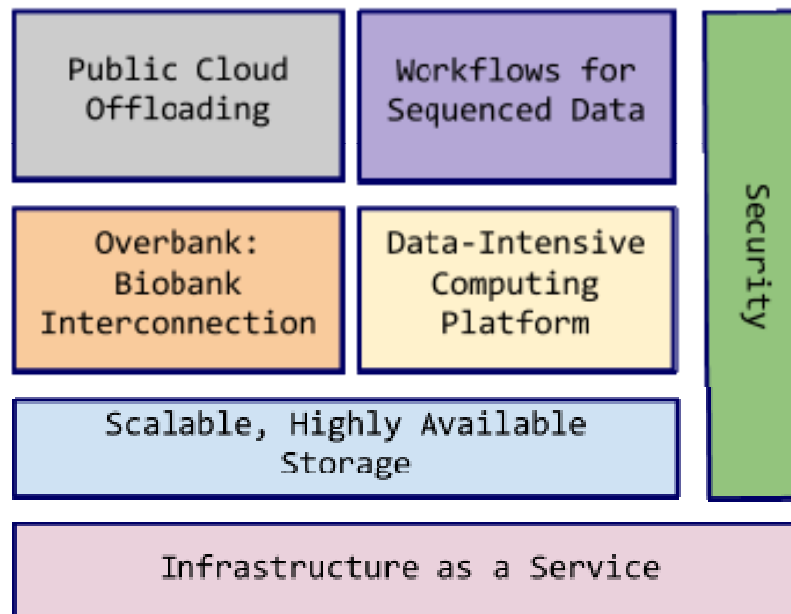
Mutable
Meta-Data



Immutable
Genomic Data

BiobankCloud PaaS

BiobankCloud PaaS



But we can't let the data leave our site...

- You buy the NGS machines
- You might want to use affordable compute and storage
-but you have to keep your genomic data on-site.



[Taken from Illumina BaseSpace]

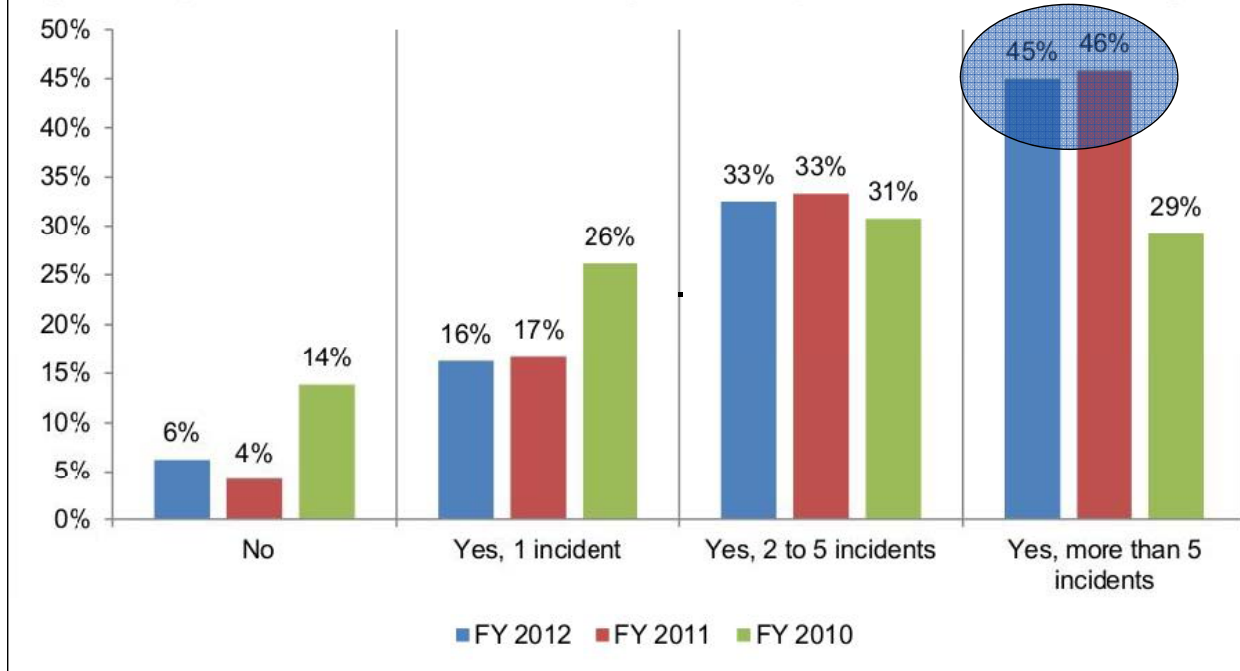
What is the real privacy?

The slide compares 'Private' and 'Public' genomic data attributes. The 'Private' side (red header) lists: Donor name, Identity number, Social security number, and Home address. The 'Public' side (green header) lists: DNA sequence, Age, and State. The background features logos for '1000 Genomes', 'Complete Genomics', 'Sanger Institute', and 'Personal Genome Project'.

Private	Public
<ul style="list-style-type: none">• Donor name• Identity number• Social security number• Home address	<ul style="list-style-type: none">• DNA sequence• Age• State

Data Leakage

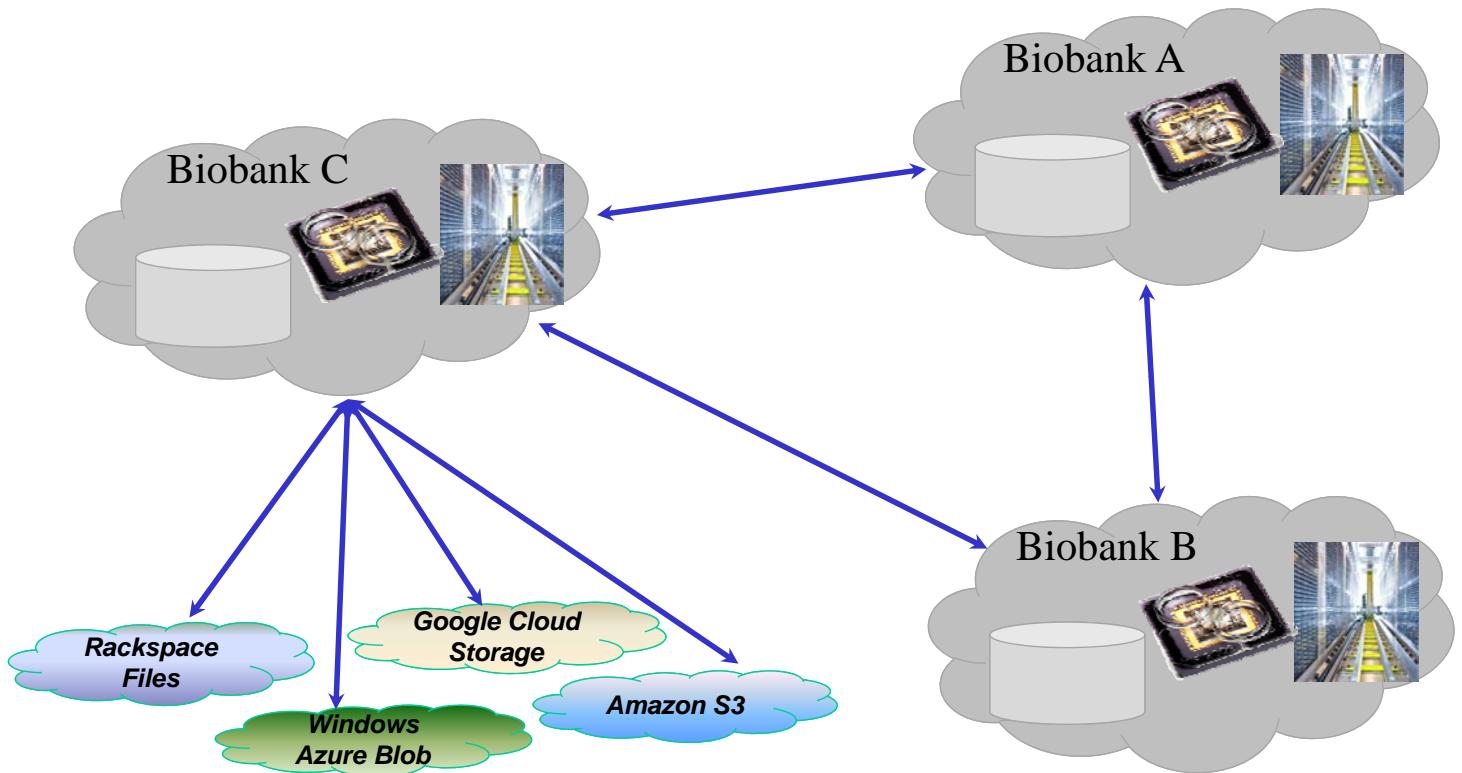
Figure 1. Experienced a data breach involving the loss of patient data in the past two years



Ponemon Institute 3rd Annual benchmark study on patient privacy and data security - 80 health care organisations

Proposed solutions

Overbank

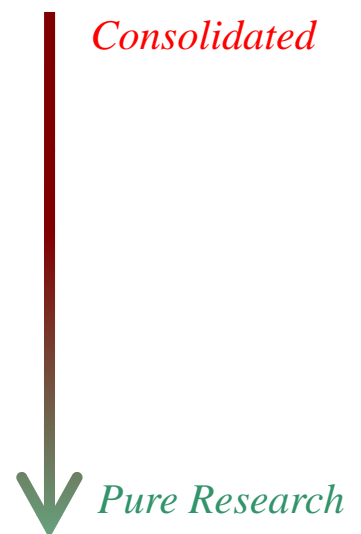


Dec. 3-4, 2012

23

Six Research Ideas

- Wide-area state machine rep.
- Cloud-of-clouds storage
- Multi-Cloud MapReduce
- Network coding
- Software-Defined Networking
- Secure two-party computation

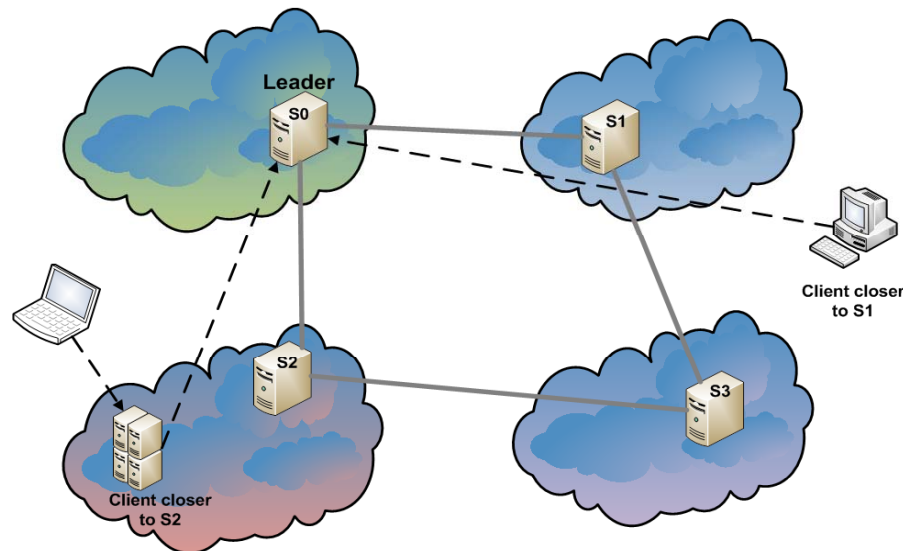


Dec. 3-4, 2012

24

Wide-Area Replication

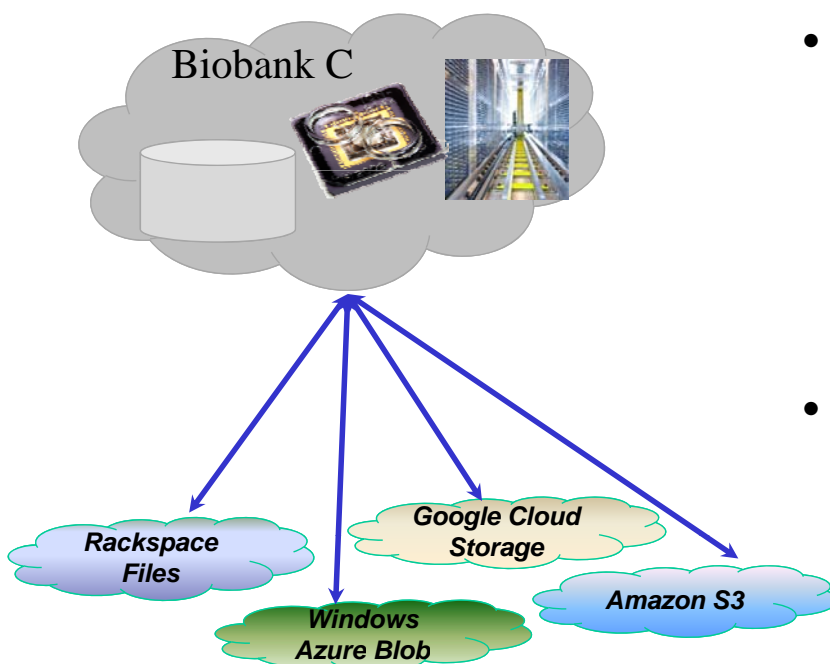
- Consistent and secure replication across datacenters



Dec. 3-4, 2012

25

Cloud-of-Clouds Storage



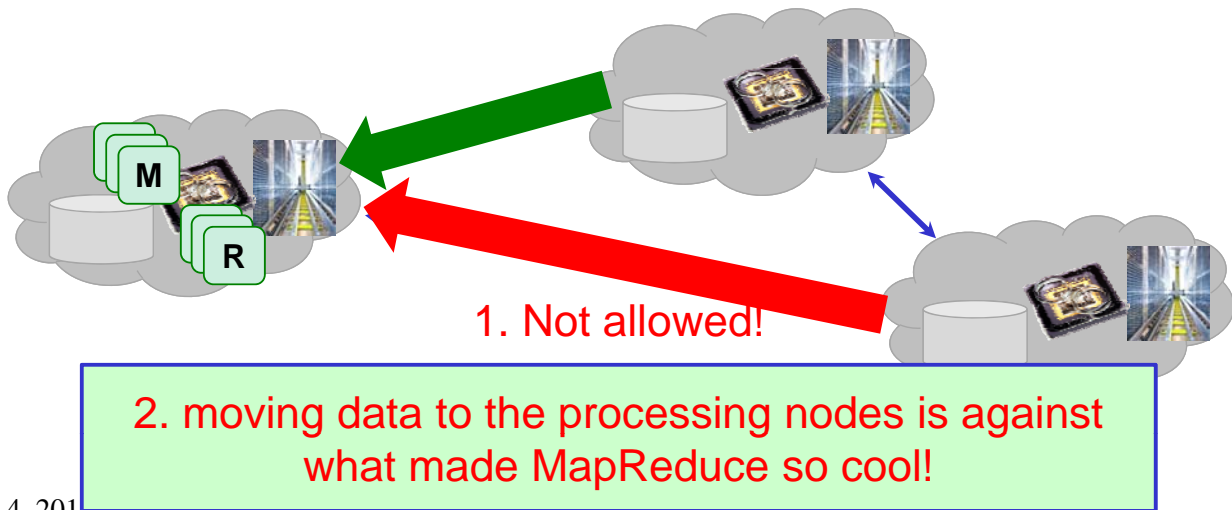
- Not everything can go to public clouds
 - Even if we encrypt and disperse the data, there are legal limitations
 - But a lot of data can!
- TClouds work: DepSky, **C2FS (file system)**
 - Can be improved for our workloads
 - Database-like system

Dec. 3-4, 2012

26

Multi-Cloud MapReduce

- What if you want to do some analysis but you need data from other Biobanks?

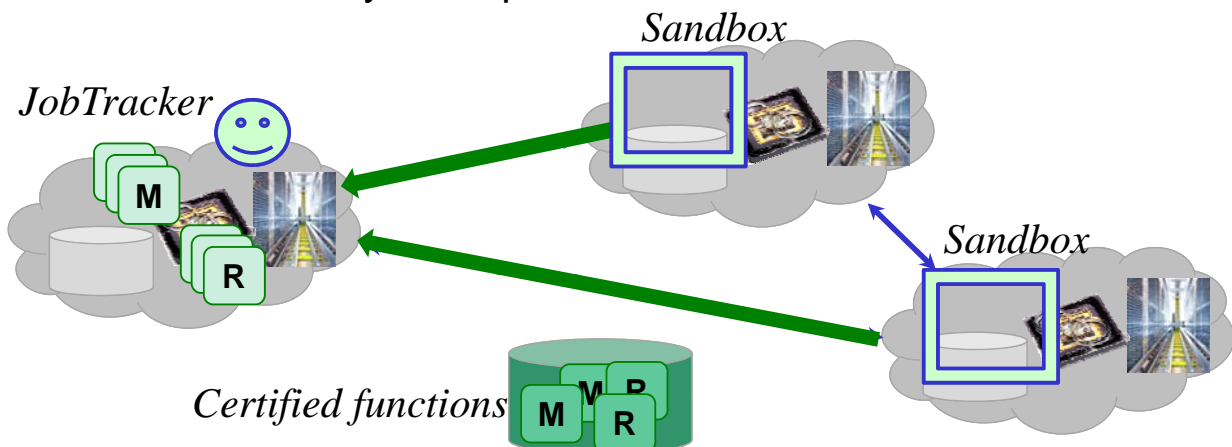


Dec. 3-4, 2012

27

Multi-Cloud MapReduce

- Move computation, not data, without privacy violations
 - sandboxing
 - wide-area job tracker
 - certified library of map and reduce functions



Dec. 3-4, 2012

28

Thank you! Tack! Danke! Obrigado!

