

# **X-Informatics**

# **Cloud Computing Technology Part II**

June 20 2013

Geoffrey Fox

[gcf@indiana.edu](mailto:gcf@indiana.edu)

<http://www.infomall.org/X-InformaticsSpring2013/index.html>

Associate Dean for Research, School of Informatics and  
Computing

Indiana University Bloomington  
2013

# **Big Data Ecosystem in One Sentence**

Use **Clouds** running **Data Analytics** processing **Big Data** to solve problems in **X-Informatics** ( or **e-X**)

X = Astronomy, Biology, Biomedicine, Business, Chemistry, Crisis, Earth Science, Energy, Environment, Finance, Health, Intelligence, Lifestyle, Marketing, Medicine, Pathology, Policy, Radar, Security, Sensor, Social, Sustainability, Wealth and Wellness with more fields (physics) defined implicitly

Spans Industry and Science (research)

Education: **Data Science** see recent New York Times articles

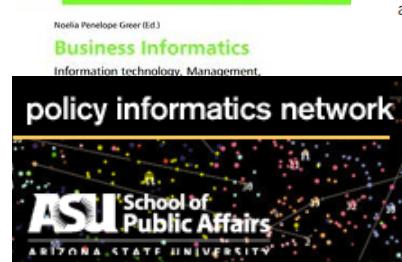
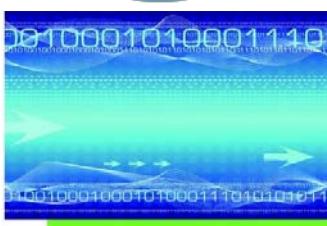
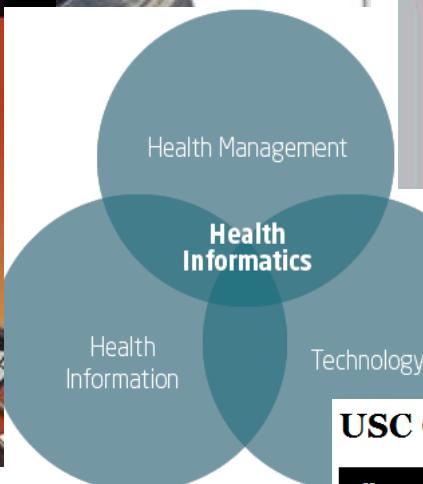
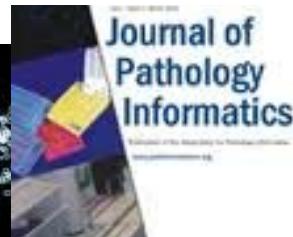
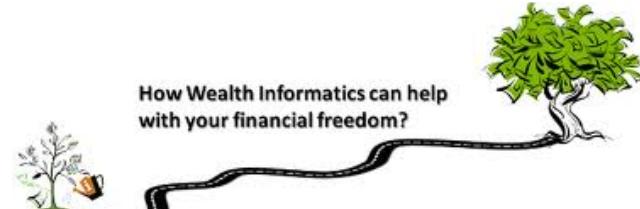
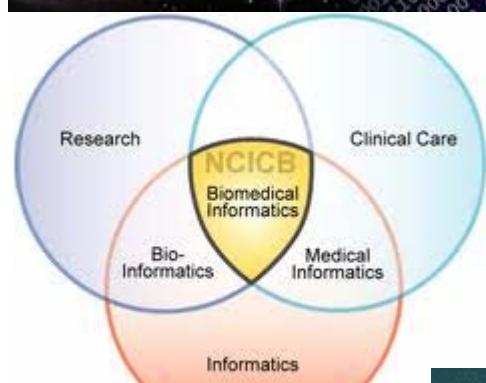
<http://datascience101.wordpress.com/2013/04/13/new-york-times-data-science-articles/>



**Earth Science  
INFORMATICS**

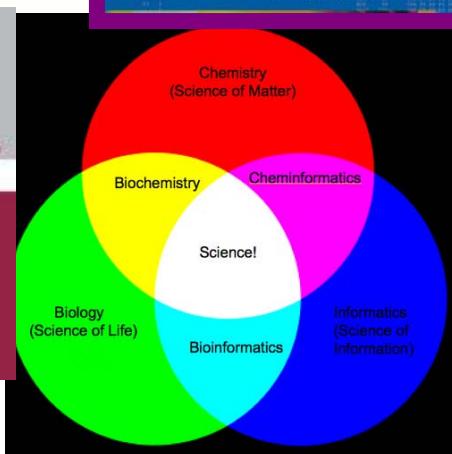
# AstroInformatics2012

Redmond, WA, September 10 - 14, 2012



# Xinformatics

**Biomedical Informatics**  
Computer Applications in Health Care  
and Biomedicine



Opportunities and Challenges  
in Crisis Informatics

## USC Center For Energy Informatics

Home Research Publications Smart



### About the Center

Welcome to the Center For Energy Informatics (CEI) at USC, an Organized Research Unit (ORU) housed in the [Viterbi School of Engineering](#). Energy Informatics is the application of info

#### Lifestyle Informatics



Applications of LI  
How is the training classified  
Occupation Prof  
Further study  
Student at the  
Watch the movie  
Studying Abro



Lifestyle Informatics: Let people live  
The study Lifestyle Informatics is about so  
this bachelor including applied psycholog  
knowledge about language and information  
short better. Lifestyle Informatics: let peo  
[Lifestyle Informatics](#)

# **What is Cloud Computing**

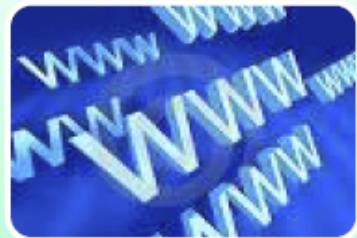
In more detail

## Cloud Computing provides solutions to a variety of challenges and opportunities



### The classical problem

- Under-utilized server resources waste computing power (and energy)
- Over-utilized servers cause interruption or degradation of service levels



### ...today in an Internet setting

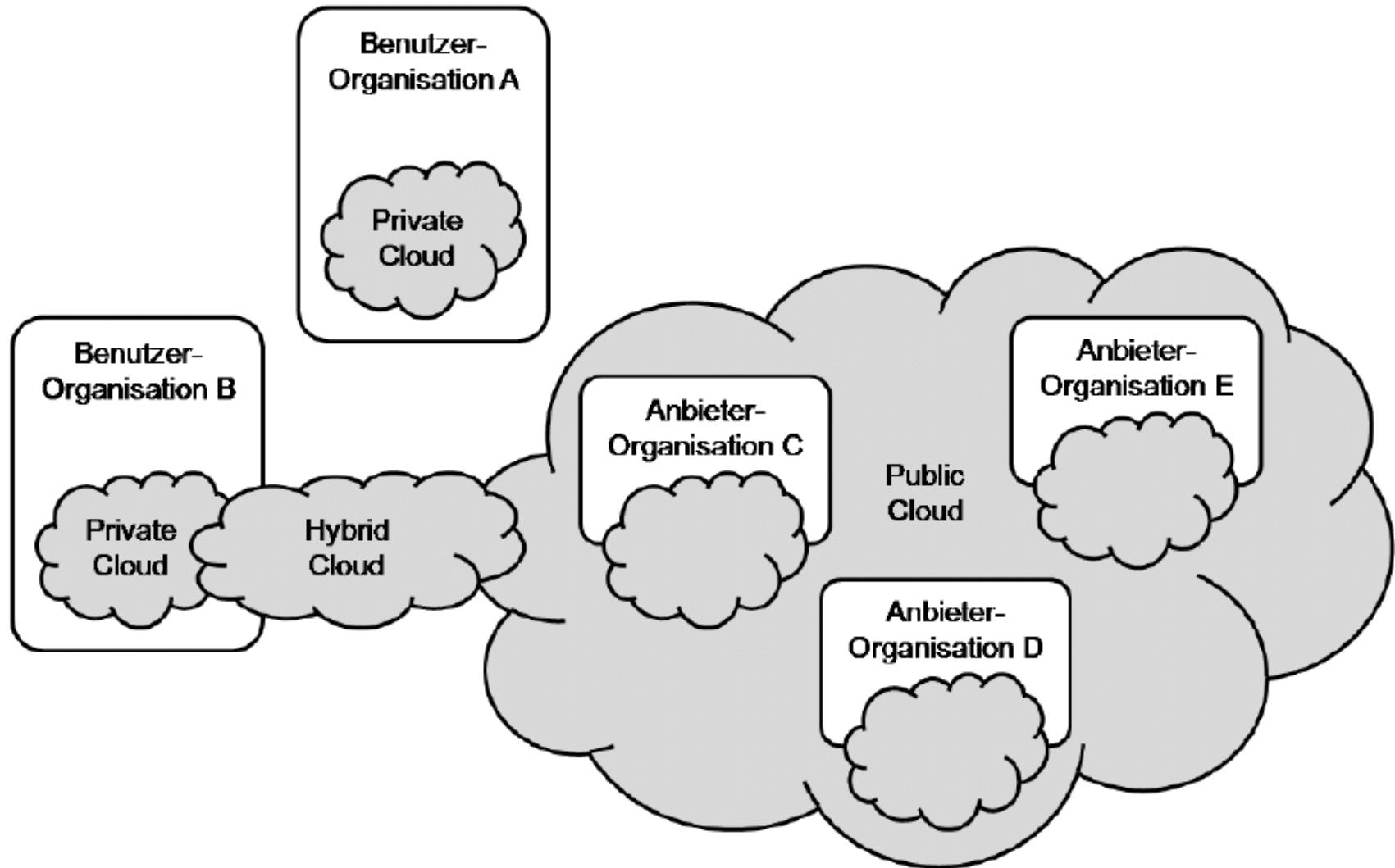
- Resource demands are increasingly of highly dynamic nature and Internet-scale
- On-demand resources are a means for faster time-to-market, and cost-effective innovation processes



### ...and tomorrow in the next-gen Web

- Leveraging the Web as a combined technology, business, and people collaboration platform:
  - Making effective use of sophisticated infrastructure which is increasingly available as (Web) services
  - Enabling dynamic (trans-)formation of open service and business networks

## Organizational Cloud Architecture: Public-/Hybrid-/Private-Cloud



## Players

Cloud **infrastructure service providers** – raw cloud resources

IaaS (infrastructure-as-a-service)

Cloud **platform providers** – resources + frameworks; PaaS (platform-as-a-service)

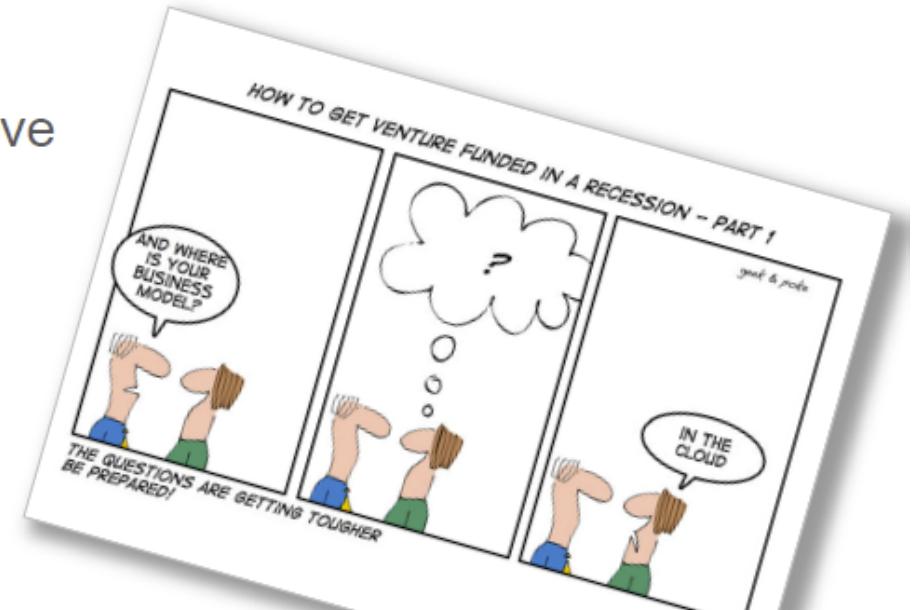
Cloud **intermediares** – help broker some aspect of raw resources and frameworks, e.g.,

server managers, application assemblers, application hosting

Cloud **application providers** (SaaS)

Cloud **consumers** – users of the above

[MM]



# Cloud Talk

- ▲ *Big “3” of Cloud Computing?*
- ▲ *Cost*
  - ▲ *Clouds are renowned for being dirt cheap for storage and burst-y processing.*
- ▲ *Flexibility*
  - ▲ *let someone else manage it for you.*
- ▲ *Elasticity*
  - ▲ *Growth and shrinkage*

johnmwillis.com



# Cloud Talk

## ▲ *More on why Clouds?*

- ▲ *On Demand Business*
  - ▲ *Unexpected loads, The <> effect*
- ▲ *Meeting Batch Load Demands*
  - ▲ *Batch*
- ▲ *Parallelism*
  - ▲ *Large clusters of parallel jobs*
- ▲ *Season Workloads*
  - ▲ *Retail, Travel, Financial*
- ▲ *Backup Storage*

johnmwillis.com



# Cloud Talk

- ▲ *What is my definition of a “Holy-Grail” Cloud*
- ▲ Abstraction of the hardware infrastructure from the service.
- ▲ Abstraction of the software infrastructure from the service.

johnmwillis.com



# Cloud Talk

## ▲ *Cloud “Services” in Simple Terms*

### ▲ *Applications*

- ▲ *Software provided as a service*

### ▲ *Middle-ware*

- ▲ *Software stacks for developers LAMP, Java application servers, .Net*

### ▲ *Servers*

- ▲ *Hardware and an operating system, perhaps a hyper-visior*

johnmwillis.com



# Cloud Talk

## ▲ *Processing Large Datasets*

### ▲ *Map Reduce*

- ▲ *Jobs that run as hundreds or even thousands of separate parallel processes.*
- ▲ *Like counting the words in a book and break it up into multiple running parts (i.e., The Map)*
- ▲ *Then collect them all back into summary counts (i.e., The Reduce.)*

johnmwillis.com



# Cloud Talk

## ▲ *Processing Large Datasets*

### ▲ *Hadoop*

- ▲ *Google invented GFS*
- ▲ *Yahoo*
- ▲ *AOL*
- ▲ *IBM*
- ▲ *Facebook*
- ▲ *Last.fm*

johnmwillis.com

<http://www.slideshare.net/botchagalupe/introduction-to-clouds-cloud-camp-columbus>



# Cloud Talk

- ▲ *Cloud Computing Challenges*
  - ▲ *Retraining developers and operations people to deal with cloud computing*
  - ▲ *Orchestration of multiple clouds*
  - ▲ *24 by 7 by 365 operations in the cloud is usually more expensive*
  - ▲ *Legacy applications might not port easily*
  - ▲ *Virtualization project disruption*
  - ▲ *Recent McKinsey Report \$366 vs \$150*



johnmwillis.com

# Cloud Talk

## ▲ *More Challenges*

- ▲ *Workload Affinities*
- ▲ *Standards (Lock-in)*
- ▲ *Weak SLA's compared to Corp*
- ▲ *Service Management*
- ▲ *Security*
- ▲ *Compliance*
- ▲ *Image Sprawl*
- ▲ *Trojan Virtual Images*
- ▲ *Governance*



johnmwillis.com

# Cloud Talk

## ▲ *New Cloud Terms*

- ▲ *Cloud Bursting*
  - ▲ *Analytics, Coding*
- ▲ *Hybrid Clouds*
  - ▲ *VPN, Multiple Clouds*
- ▲ *Cloud Spillage*
  - ▲ *An IBM Term*
- ▲ *Cloud Orchestration*
  - ▲ *Managing multiple clouds*

johnmwillis.com



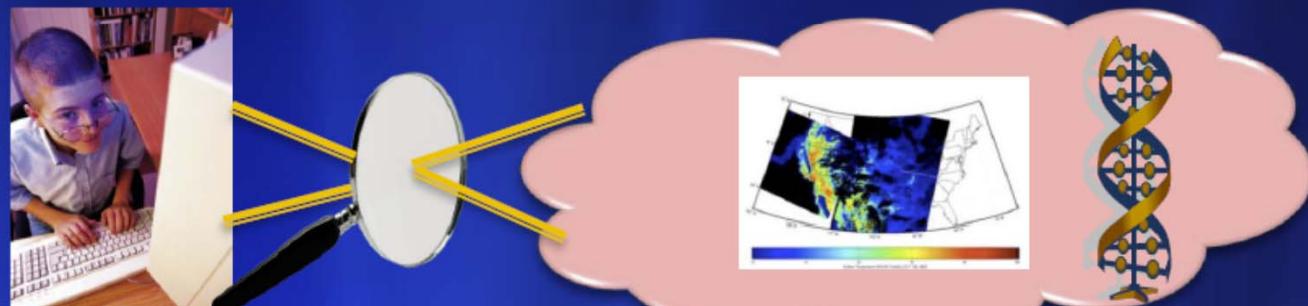
# The Cloud as an extension of your desktop and other client devices

- Today

- Cloud storage for your data files synchronized across all your machines (mobile me, live mesh, flicker, etc.)
- Your collaboration space (Sakai, SharePoint)
- Cloud-enabled apps (Google Apps, Office Live)

- Tomorrow (or even sooner)

- The lens that magnifies the power of desktop
- Operate on a table with a billion rows in excel
- Matlab analysis of a thousand images in parallel



# Cloud Properties

- **Designed to Provide Information and Computation to Many Users**
- Automatic Deployment and Management of Virtual Machine Instances
  - tens to thousands & dynamic scalability
- Dynamic Fault Recovery of Failed Resources
  - Cloud Services must run 24x7
- Automatic Data Replication
  - Geo-replication if needed
- Two levels of parallelism
  - Thousands of concurrent users
  - Thousands of servers for a single task.



# What is the **cloud**?



Delivering Information Technology as a Standardized Service

The illusion of infinite compute and storage

Customer datacenter

Partner datacenter

Public datacenter

[http://research.microsoft.com/en-us/um/redmond/events/cloudfutures2012/tuesday/Keynote\\_OpportunitiesAndChallenges\\_Yousef\\_Khalidi.pdf](http://research.microsoft.com/en-us/um/redmond/events/cloudfutures2012/tuesday/Keynote_OpportunitiesAndChallenges_Yousef_Khalidi.pdf)

# Clouds Offer From different points of view

- **Features from NIST:**
  - On-demand service (elastic);
  - Broad network access;
  - Resource pooling;
  - Flexible resource allocation;
  - Measured service
- **Economies of scale** in performance and electrical power (**Green IT**)
- Powerful new **software models**
  - **Platform as a Service** is not an alternative to **Infrastructure as a Service** – it is instead an incredible valued added
  - Amazon is as much PaaS as Azure

# What is Cloud?

- **Defining the Cloud**
  - A scalable, persistent outsourced infrastructure
  - A framework for massive data analysis
  - An amplifier of our desktop experience
- **The Origins**
  - Modern data center architecture

[http://research.microsoft.com/en-us/people/barga/sc09\\_cloudcomp\\_tutorial.pdf](http://research.microsoft.com/en-us/people/barga/sc09_cloudcomp_tutorial.pdf)

# What is a cloud?



SLAs

The Cloud is

data center architecture

Web Services



Virtualization

# Defining the Cloud

- A model of computation and data storage based on “pay as you go” access to unlimited remote data center capabilities.
- A cloud infrastructure provides a framework to manage scalable, reliable, on-demand access to applications.
- Examples:
  - Search, email, social networks
  - File storage (Live Mesh, Mobile Me, Flickr, ...)
- A way for a start-up to build a scalable web presence without purchasing hardware.



[http://research.microsoft.com/en-us/people/barga/sc09\\_cloucomp\\_tutorial.pdf](http://research.microsoft.com/en-us/people/barga/sc09_cloucomp_tutorial.pdf)

# Cloud Mythologies

- Cloud computing infrastructure is just a web service interface to operating system virtualization.
  - “I’m running Xen in my data center – I’m running a private cloud.”
- Clouds and Grids are equivalent
  - “In the mid 1990s, the term grid was coined to describe technologies that would allow consumers to obtain computing power on demand.”
- Cloud computing imposes a significant performance penalty over “bare metal” provisioning.
  - “I won’t be able to run a private cloud because my users will not tolerate the performance hit.”

# The Clients+Cloud Platform

- At one time the “client” was a PC + browser.
- Now the cloud is an *integration point* for
  - The Phone
  - The laptop/tablet
  - The TV/Surface/Media wall
- And the future
  - The instrumented room
  - Aware and active surfaces
  - Voice and gesture recognition
  - Knowledge of where we are
  - Knowledge of our health



# The History of the Cloud

- In the beginning ...
  - There was search, email, messaging, web hosting
- The challenge: How do you
  - Support email for 375 million users?
  - Store and index 6.75 trillion photos?
  - Support 10 billion web search queries/month?
  - Build an index for the entire web? And do it over and over again...
- And
  - deliver a quality response in 0.15 seconds to millions of simultaneous users?
  - never go down.
- Solution: build big data centers

# Public, Private, and Premise



- **Public Cloud**

- Large scale infrastructure available on a rental basis
- Virtualized compute, network and storage
- Underlying infrastructure is shared but tenants are isolated
- Interface is transactional
- Accounting is e-commerce based

- **Private Cloud**

- Dedicated resources either as a rental or on-premise

- **On-premise Cloud**

- Like public clouds but
  - Isolation must be controllable
  - Accounting is organizational

# Public IaaS



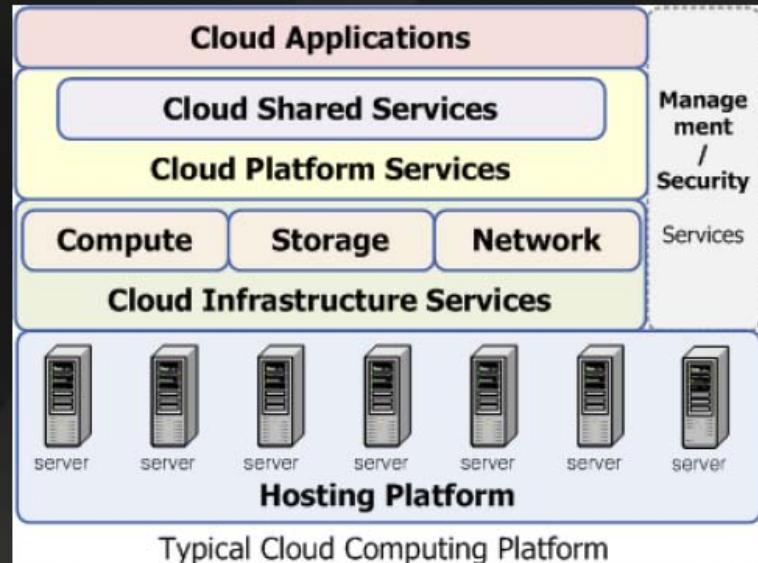
- **Large scale infrastructure available on a rental basis**
  - Operating System virtualization (e.g. Xen, KVM) provides CPU isolation
  - “Roll-your-own” network provisioning provides network isolation
  - Locally specific storage abstractions
- **Fully customer self-service**
  - Customer-facing Service Level Agreements (SLAs) are advertized
  - Requests are accepted and resources granted via web services
  - Customers access resources remotely via the Internet
- **Accountability is e-commerce based**
  - Web-based transaction
  - “Pay-as-you-go” and flat-rate subscription
  - Customer service, refunds, etc.

# **As a Service and Platform Model**

IaaS PaaS SaaS

# Typical Cloud Computing Platform

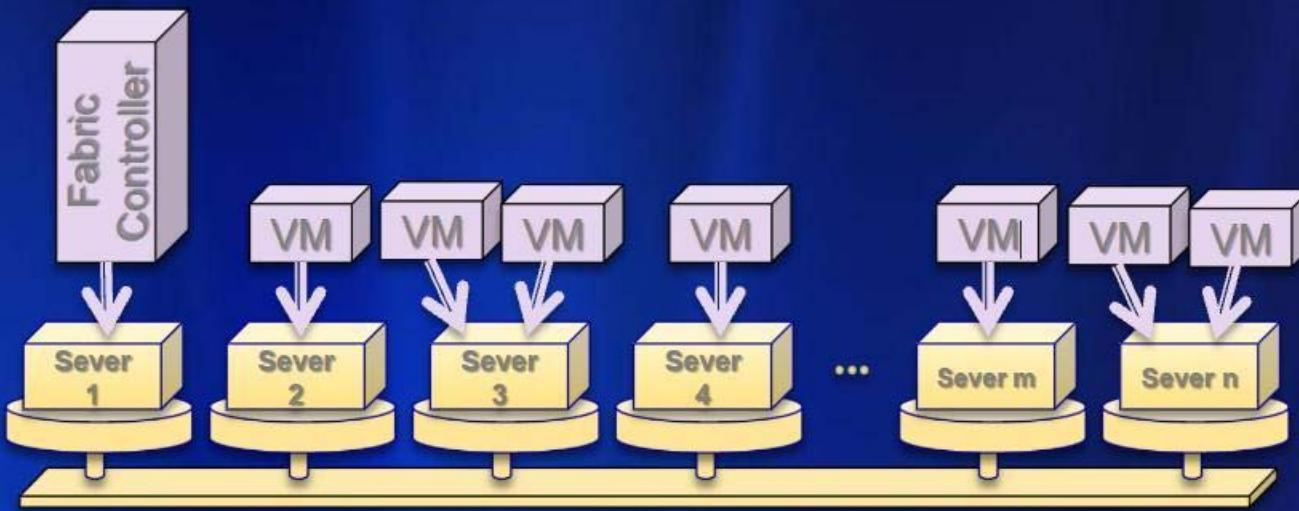
- Hosting Platform
  - Provides the physical, virtual, and software assets which include physical machines, operating systems, network systems, storage systems, power management, and virtualization software
- Cloud Infrastructure Services(IaaS)
  - Abstract the hosting platform as a set of virtual resources(i.e. compute, storage, and network)
  - Manage those resources based on scalability and availability needs
- Cloud Platform Services(PaaS)
  - Provide a set of capabilities exposed as a services to help with integrating on-premise software with hosted services
- Cloud Applications(SaaS)
  - Houses applications that are built for cloud computing, which expose Web interfaces and Web Services for end users, enabling multitenant hosting models.



- Security Services
  - Ensure token provisioning, identity federation, and claims transformation
- Management Services
  - Provide a set of capabilities to automate scalability and availability administration such as deployment configurations, service usage analytics, and connection to enterprise management systems

# Three Levels of Cloud Architecture

- Infrastructure as a Service (IaaS)
  - Provide App builders a way to configure a Virtual Machine and deploy one or more instances on the data center
  - Each VM has access to local and shared data storage
  - The VM has an IP Address visible to the world
  - A Fabric controller manages VM instances
    - Failure and restart, dynamic scale out and scale back.



# Windows Azure – a platform for apps



Physical



Virtual



IaaS



PaaS



SaaS

The Foundation for  
Private Cloud

You can build apps that use both

Finished Services

[http://research.microsoft.com/en-us/um/redmond/events/cloudfutures2012/tuesday/Keynote\\_OpportunitiesAndChallenges\\_Yousef\\_Khalidi.pdf](http://research.microsoft.com/en-us/um/redmond/events/cloudfutures2012/tuesday/Keynote_OpportunitiesAndChallenges_Yousef_Khalidi.pdf)

# Clouds have highlighted SaaS PaaS IaaS

**Software  
(Application  
Or Usage)**

**SaaS**

- Education
- Applications
- CS Research Use e.g. test new compiler or storage model

**Platform  
PaaS**

- Cloud e.g. MapReduce
- HPC e.g. PETSc, SAGA
- Computer Science e.g. Compiler tools, Sensor nets, Monitors

**Infra  
structure**

**IaaS**

- Software Defined Computing (virtual Clusters)
- Hypervisor, Bare Metal
- Operating System

**Network  
NaaS**

- Software Defined Networks
- OpenFlow GENI

**But equally valid for classic clusters**

- Software Services are building blocks of applications
- The middleware or computing environment including **HPC, Grids ...**
- Nimbus, Eucalyptus, OpenStack, OpenNebula CloudStack plus **Bare-metal**
- OpenFlow – *likely to grow in importance*

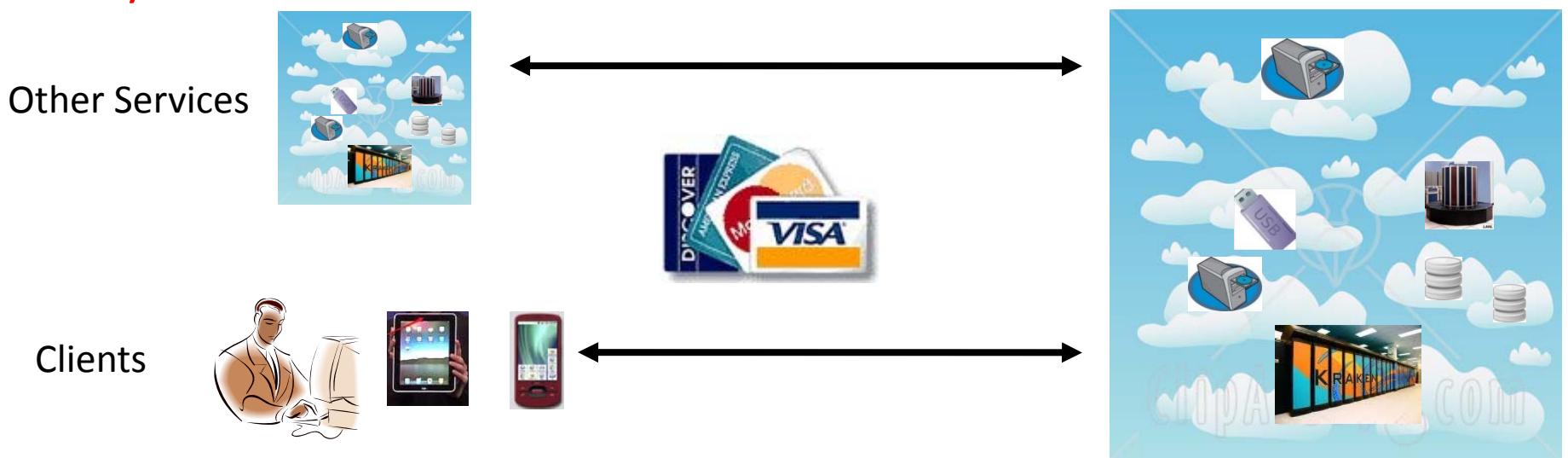


# IaaS Open Source

- **Eucalyptus** (first major open source VM management environment)
- **Nimbus**
- **OpenStack** (currently dominant – 2800 at last OpenStack meeting)
- **OpenNebula** (Europe)
- **CloudStack** (break away from OpenStack)
- These manage basic computing and storage but these are separable ideas
- Management, Networking, Access Tools
  
- You install on your own clusters and do have control over operation
- Compare Public Clouds (Amazon Azure Google) where operation is opaque but you just need a credit card and a web browser
- Compare VMware commercial version which is similar to OpenSource but with better support and probably greater robustness/functionality
- Eucalyptus has Open Source and Commercial versions

# X as a Service

- **SaaS: Software as a Service** imply software capabilities (programs) have a service (messaging) interface
  - Applying systematically reduces system complexity to being linear in number of components
  - Access via messaging rather than by installing in /usr/bin
- **IaaS: Infrastructure as a Service** or **HaaS: Hardware as a Service** – get your computer time with a credit card and with a Web interface
- **PaaS: Platform as a Service** is **IaaS** plus core software capabilities on which you build **SaaS**
- **Cyberinfrastructure** is “Research as a Service”



# PaaS – What is a "cloud platform"?



“... data as a service...”

“cloud computing journal reports that...”

“... software as a service...”

“... everything as a service...”

Platforms succeed when the platform helps others succeed

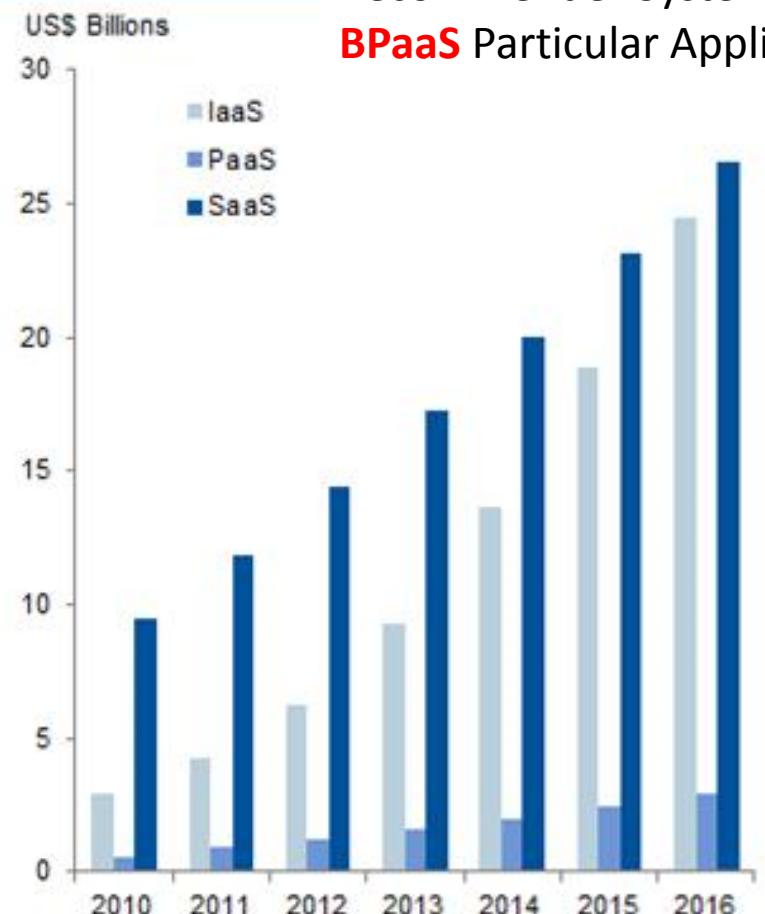
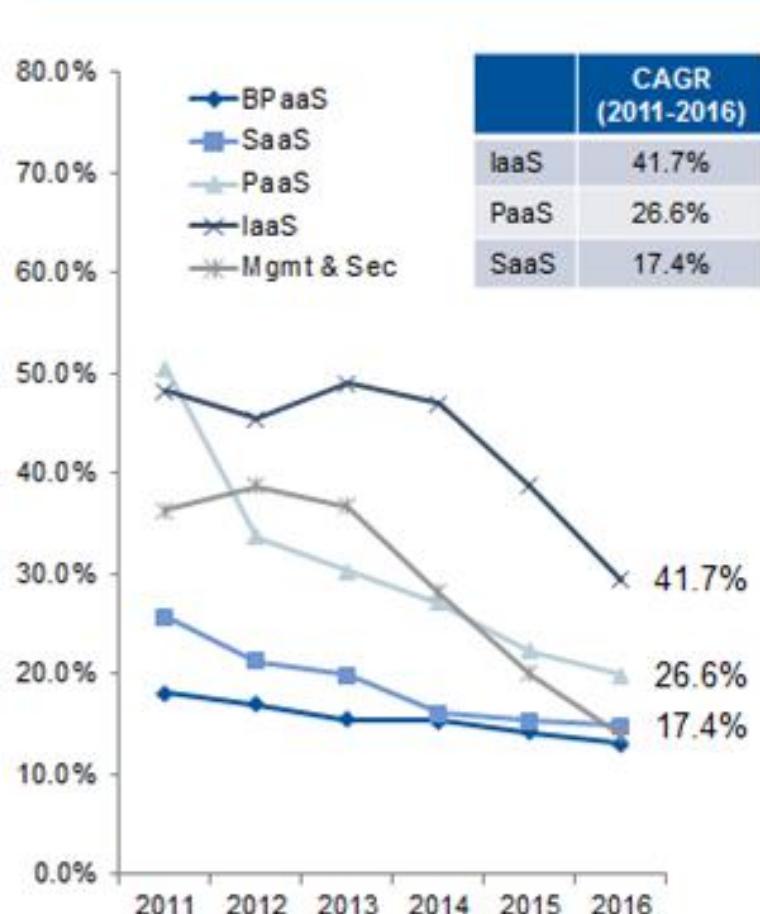
[http://research.microsoft.com/en-us/people/barga/sc09\\_cloudcomp\\_tutorial.pdf](http://research.microsoft.com/en-us/people/barga/sc09_cloudcomp_tutorial.pdf)

# Software as a Service

- Online delivery of applications
- Via Browser
  - Microsoft Office Live Workspace
  - Google Docs, etc.
  - File synchronization in the cloud – Live Mesh, Mobile Me
  - Social Networks, Photo sharing, Facebook, wikipedia etc.
- Via Rich Apps
  - Science tools with cloud back-ends
    - Matlab, Mathematica
  - Mapping
    - MS Virtual Earth, Google Earth
  - Much more to come.

# High Growth Expected in Cloud Infrastructure Services

**IaaS** Hardware e.g. Server  
**PaaS** Systems Services e.g. MapReduce, Database  
**SaaS** Applications e.g. Recommender System, Clustering  
**BPaaS** Particular Application Set



Source: Public Cloud Services Forecast, 2Q12 Update

18

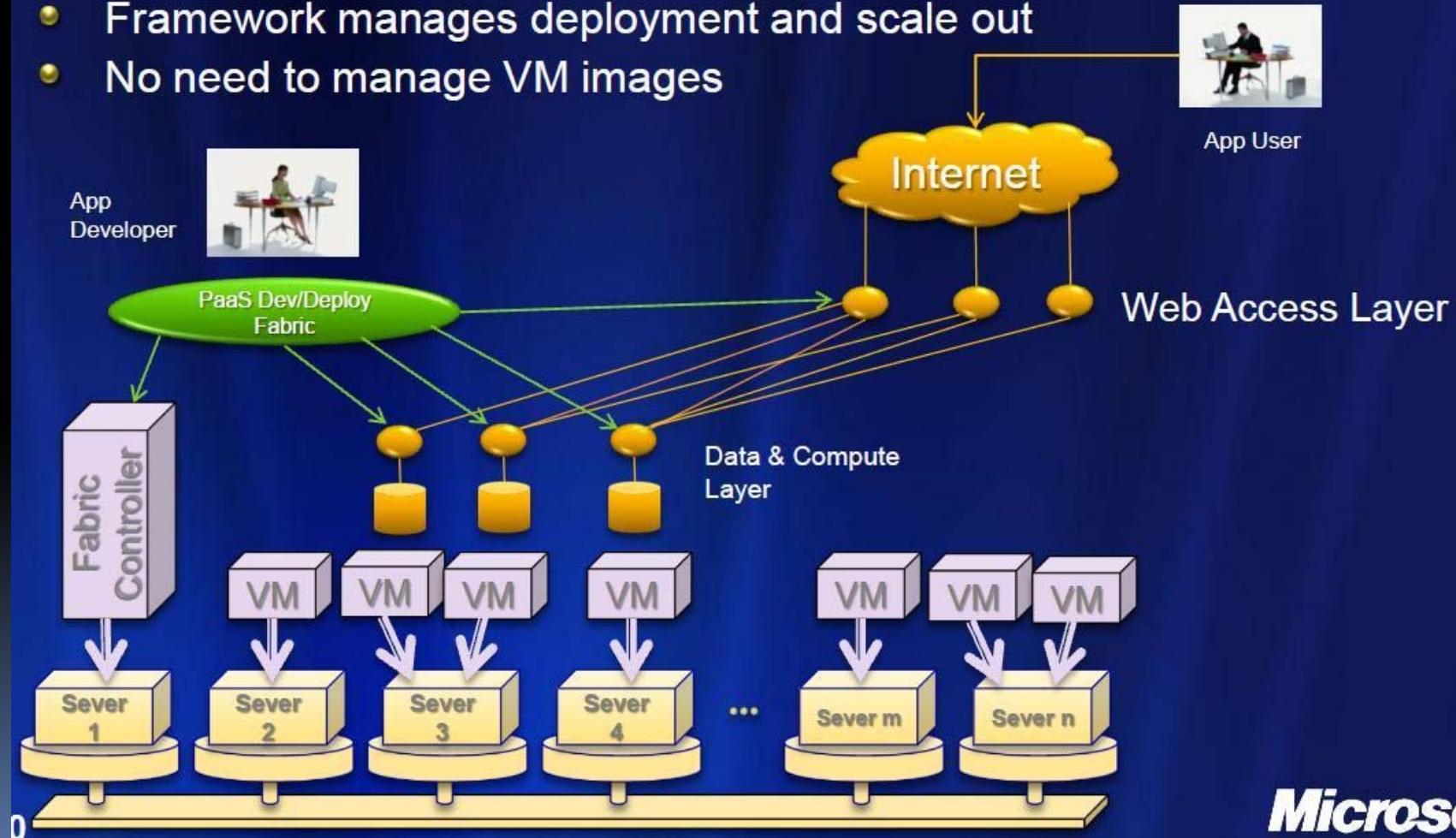
Gartner

BPM = Business Process management

# **Platform as a Service**

# Platform as a Service

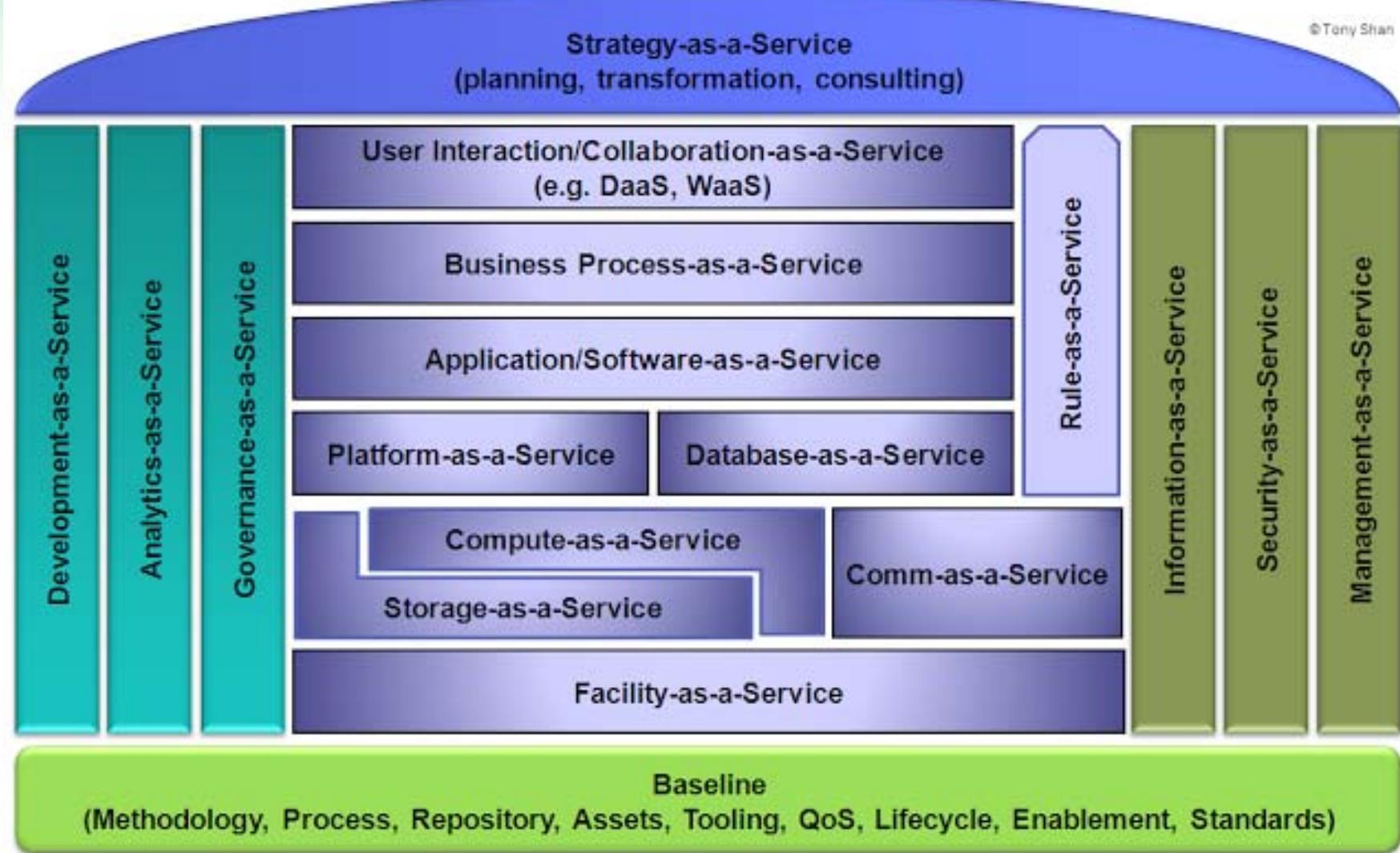
- An application development, deployment and management fabric.
- User programs web service front end and computational & Data Services
- Framework manages deployment and scale out
- No need to manage VM images



# Cloud Computing: Infrastructure and Runtimes

- **Cloud Infrastructure:** outsourcing of servers, computing, data, file space, utility computing, etc.
  - Handled through (Web) services that control virtual machine lifecycles.
- **Cloud Runtimes or Platform:** tools (for using clouds) to do data-parallel (and other) computations.
  - Apache Hadoop, Google MapReduce, Microsoft Dryad, Bigtable, Chubby (synchronization) and others
  - MapReduce designed for information retrieval but is excellent for a wide range of **science data analysis applications**
  - Can also do much traditional parallel computing for data-mining if extended to support **iterative** operations
  - MapReduce not usually done on Virtual Machines

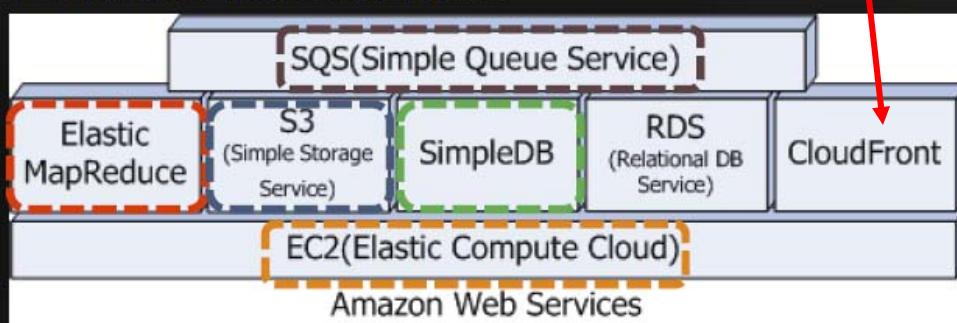
Grid/Cloud	<b>Authentication and Authorization:</b> Provide single sign in to both FutureGrid and Commercial Clouds linked by workflow
	<b>Workflow:</b> Support workflows that link job components between FutureGrid and Commercial Clouds. Trident from Microsoft Research is initial candidate
	<b>Data Transport:</b> Transport data between job components on FutureGrid and Commercial Clouds respecting custom storage patterns
	<b>Software as a Service:</b> This concept is shared between Clouds and Grids and can be supported without special attention
	<b>SQL:</b> Relational Database
Cloud	<b>Program Library:</b> Store Images and other Program material (basic FutureGrid facility)
	<b>Blob:</b> Basic storage concept similar to Azure Blob or Amazon S3
	<b>DPFS Data Parallel File System:</b> Support of file systems like Google (MapReduce), HDFS (Hadoop) or Cosmos (Dryad) with compute-data affinity optimized for data processing
	<b>Table:</b> Support of Table Data structures modeled on Apache Hbase (Google Bigtable) or Amazon SimpleDB/Azure Table (eg. Scalable distributed "Excel")
	<b>Queues:</b> Publish Subscribe based queuing system
	<b>Worker Role:</b> This concept is implicitly used in both Amazon and TeraGrid but was first introduced as a high level construct by Azure
	<b>Web Role:</b> This is used in Azure to describe important link to user and can be supported in FutureGrid with a Portal framework
	<b>MapReduce:</b> Support MapReduce Programming model including Hadoop on Linux, Dryad on Windows HPCS and Twister on Windows and Linux



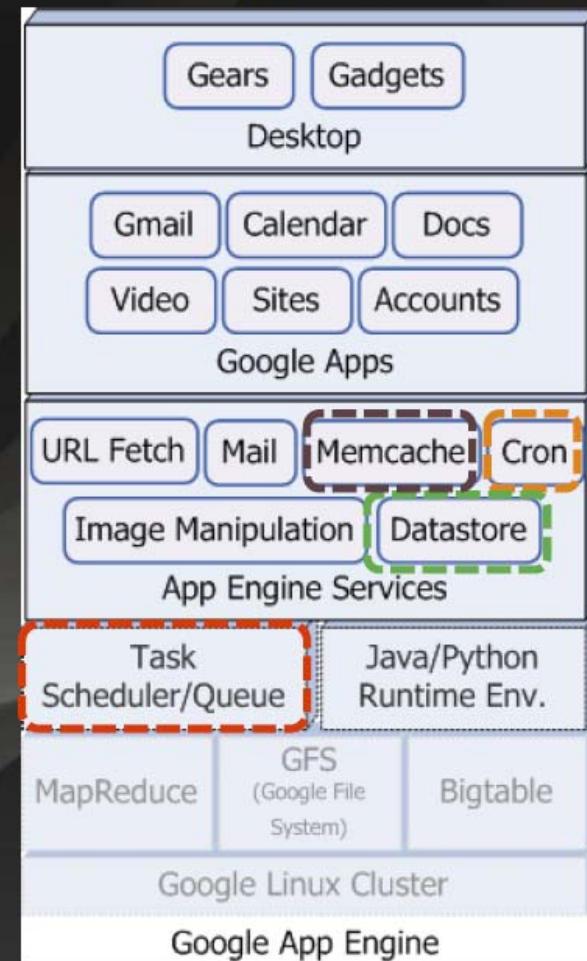
<http://cloudonomic.blogspot.com/2009/02/cloud-taxonomy-and-ontology.html>

# Architectures of Public Cloud Computing

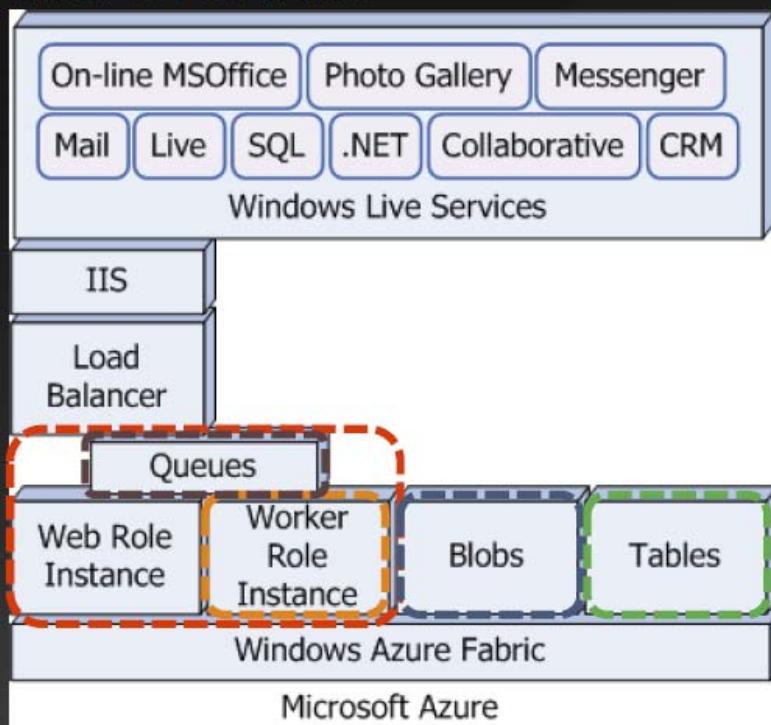
- Amazon Web Services



- Google App Engine



- Microsoft Azure



# Platform Extension to Cloud is a Continuum

Less Constrained

## Constraints in the App Model

More Constrained

### Amazon AWS

VMs Look Like Hardware  
No Limit on App Model  
User Must Implement Scalability and Failover

### Microsoft Azure

.NET CLR/Windows Only  
Choice of Language  
Some Auto Failover/  
Scale (but needs declarative application properties)

### Google App Engine

Traditional Web Apps  
Auto Scaling and Provisioning

### Force.Com

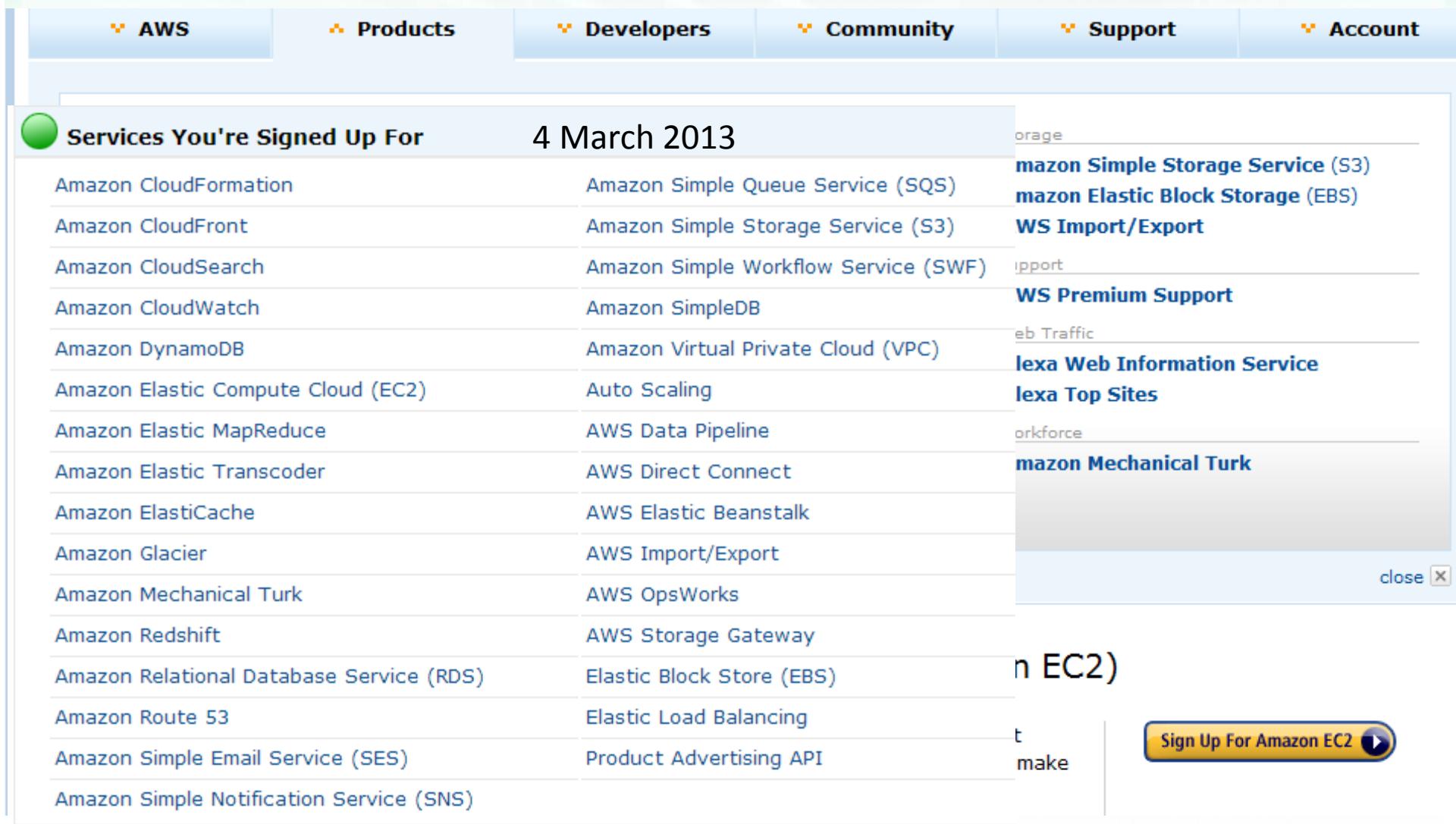
SalesForce Biz Apps  
Auto Scaling and Provisioning

Less Automation

## Automated Management Services

More Automation

# Amazon offers a lot! i.e. is PaaS



The screenshot shows a metrics dashboard for an Amazon Lambda function named "lambdafunction1". The left sidebar lists metrics: CPU Usage, Network In, Network Out, and Memory Utilization. The main area displays a line chart of CPU Usage over a 1-hour period, with values ranging from 0% to 100%. A tooltip for the peak value shows 100% at 2013-03-04T12:00:00Z.

**Signed Up Services**

Service	Description
Amazon CloudFormation	CloudFormation makes it easy to create and manage highly reliable, highly available, and completely portable AWS infrastructure with just a few clicks or a simple command line interface.
Amazon CloudFront	CloudFront is a global content delivery network that makes it easy to serve content from multiple locations.
Amazon CloudSearch	CloudSearch makes it easy to search large amounts of unstructured text data.
Amazon CloudWatch	CloudWatch provides monitoring and management for your Amazon Web Services (AWS) environment.
Amazon DynamoDB	DynamoDB is a fast, fully managed NoSQL database service.
Amazon Elastic Compute Cloud (EC2)	EC2 provides secure, resizable compute capacity in the cloud.
Amazon Elastic MapReduce	Elastic MapReduce makes it easy to process large amounts of data using MapReduce.
Amazon Elastic Transcoder	Elastic Transcoder makes it easy to convert video files between formats.
Amazon ElastiCache	ElastiCache makes it easy to store and retrieve small amounts of data at high rates.
Amazon Glacier	Glacier is a low-cost storage service for infrequently accessed data.
Amazon Mechanical Turk	Mechanical Turk makes it easy to build distributed computing applications.
Amazon Redshift	Redshift is a fast, fully managed data warehouse service.
Amazon Relational Database Service (RDS)	RDS makes it easy to set up, operate, and scale a relational database in the cloud.
Amazon Route 53	Route 53 is a highly available, distributed Domain Name System (DNS) web service.
Amazon Simple Email Service (SES)	SES makes it easy to send transactional and promotional emails.
Amazon Simple Notification Service (SNS)	SNS makes it easy to publish and subscribe to real-time notifications.
Amazon Simple Queue Service (SQS)	SQS makes it easy to build message-driven architectures.
Amazon Simple Storage Service (S3)	S3 is a highly durable and accessible cloud storage service.
Amazon Elastic Block Storage (EBS)	EBS provides persistent block storage for Amazon EC2 instances.
WS Import/Export	Import/Export makes it easy to move data between Amazon S3 and other data stores.
WS Premium Support	Premium Support provides 24x7 support for Amazon services.
Web Traffic	Web Traffic provides insights into website traffic and user behavior.
Lexa Web Information Service	Lexa Web Information Service provides access to a wide range of web-based information.
Lexa Top Sites	Lexa Top Sites provides access to the most popular websites on the web.
Mechanical Turk	Mechanical Turk provides access to a global workforce of human workers.

**Sign Up For Amazon EC2**

# **Data in the Cloud**

# Cloud Data Models

- Clouds have a unique processing model – MapReduce with Hadoop which we discussed
- They have also inspired a revolution in data models
  - Originally enterprises used databases but
  - Science did not as not efficient and “transactions” not typical use
  - Rather tended to access large chunks of data and/or did not do SQL queries. Rather used custom user software to access events of interest e.g. events like the Higgs
  - LHC explored use of “object oriented” databases (Objectivity) but abandoned
- “Web 2.0” applications had same issues as science
  - Huge data sizes (in fact larger than science) and non SQL queries
- This spurred NOSQL approaches to data and approaches to parallelism

# Different formats for data

- Traditional Windows or **UNIX file systems** arranged in a hierarchy with directories
  - Generalize to share files across many computers
  - Access with standard UNIX I/O commands and in growing fashion via web (as in Oncourse or Dropbox)
- **Google File System** – realized in Hadoop File System HDFS – splits files up for parallelism
  - Traditional file systems also had parallel extensions but GFS very explicit
- **Object stores** like Amazon S3 have simpler interface than traditional file systems
- **Databases** will appear on a later page as a higher level management system

# Clouds as Support for Data Repositories?

- The **data deluge** needs cost effective computing
  - Clouds are by definition cheapest
  - Need data and computing co-located
- **Shared resources** essential (to be cost effective and large)
  - Can't have every scientists downloading petabytes to personal cluster
- Need to reconcile **distributed** (initial source of ) **data** with shared analysis
  - Can move data to (discipline specific) clouds
  - How do you deal with multi-disciplinary studies
- **Data repositories of future will have cheap data and elastic cloud analysis support?**
  - Hosted free if data can be used commercially?

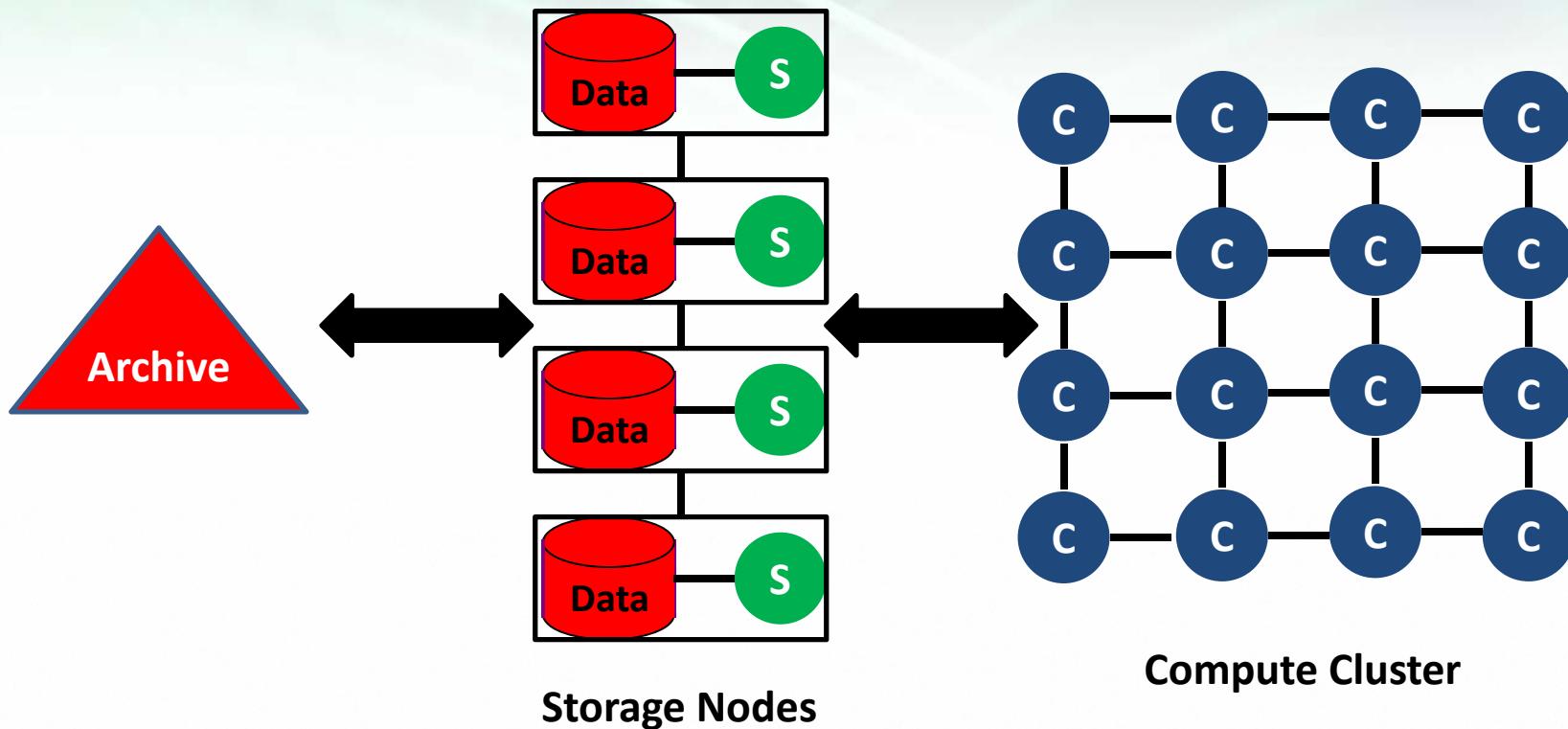
# Architecture of Data Repositories?

- Traditionally governments set up repositories for data associated with particular missions
  - For example EOSDIS (Earth Observation), GenBank (Genomics), NSIDC (Polar science), IPAC (Infrared astronomy)
  - LHC/OSG computing grids for particle physics
- This is complicated by volume of data deluge, distributed instruments as in gene sequencers (maybe centralize?) and need for intense computing like Blast
  - i.e. **repositories need lots of computing?**



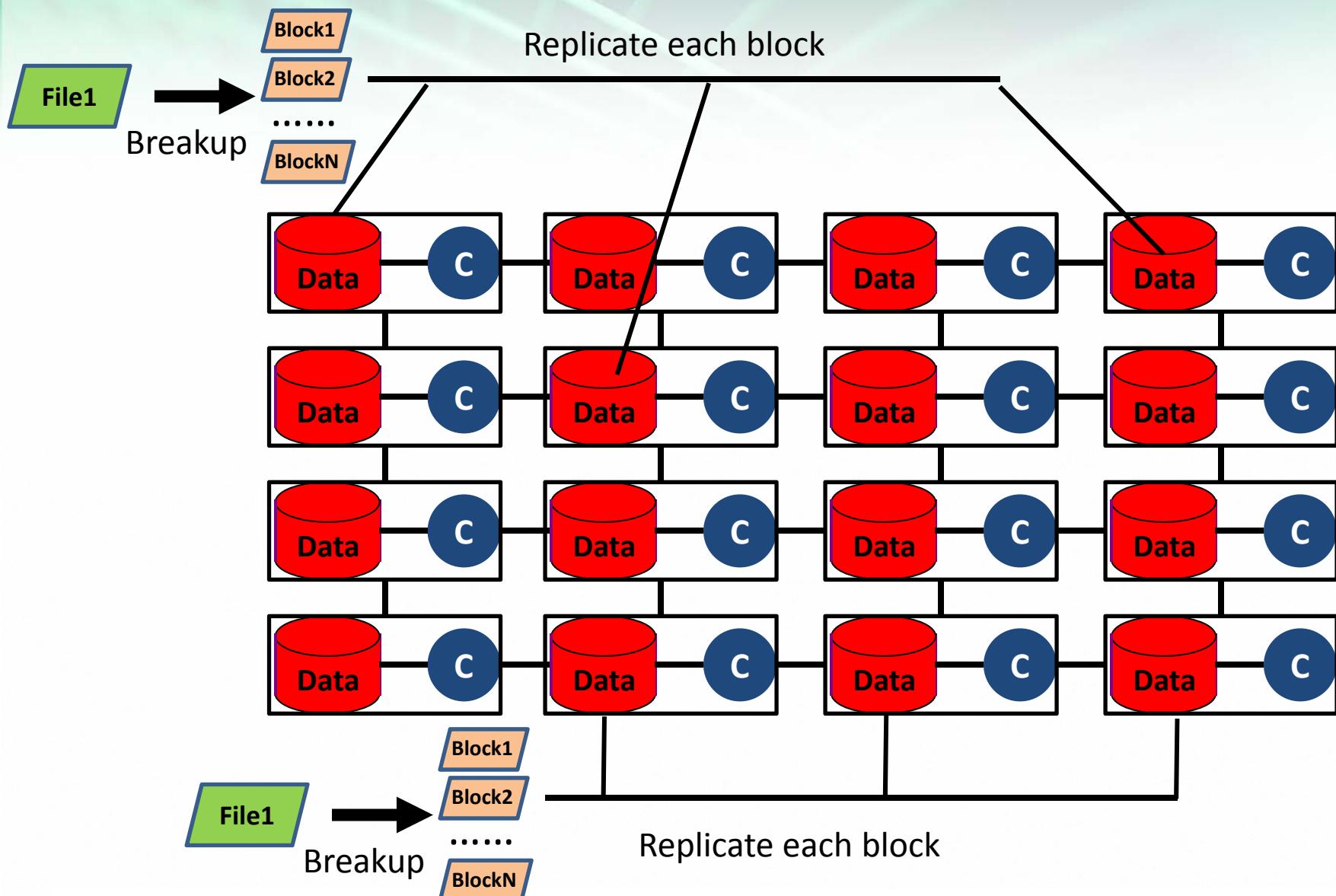
<https://portal.futuregrid.org>

# Traditional File System?



- Typically a shared file system (Lustre, NFS ...) used to support high performance computing
- Big advantages in flexible computing on shared data but doesn't **"bring computing to data"**
- Object stores similar structure (separate data and compute) to this

# Data Parallel File System?



- No archival storage and computing brought to data

# Different Approaches to managing Data

- Traditional is Directory structure and file names
  - With an add-on metadata store usually put in a database
  - Contents of file also can be self managed as with XML files
- Also traditional is databases which have tables and indices as core concept
- Newer NOSQL approaches are Hbase, Dynamo, Cassandra, MongoDB, CouchDB,Riak
  - Amazon SimpleDB and Azure table on public clouds
  - These often use an HDFS style storage and store entities in distributed scalable fashion
  - No fixed Schema
  - “Data center” or “web” scale storage
  - Document stores and Column stores depending on decomposition strategy
  - Typically linked to Hadoop for processing



<https://portal.futuregrid.org>