

25th Birthday



Cloud computing:
From technological advances
to scientific challenges

Luc Bougé
ENS Cachan/Rennes, IRISA, INRIA

With help from many colleagues:
Gabriel Antoniu, Guillaume Pierre, Louis-Claude Canon, etc.



Plan

- Introduction: where do cloud come from?
The momentum toward cloud
- Cloud computing technology: what made it possible
- Cloud computing today: For real!
- Some scientific challenges about clouds:
zooming on some recent research
- What's next? Help welcome!

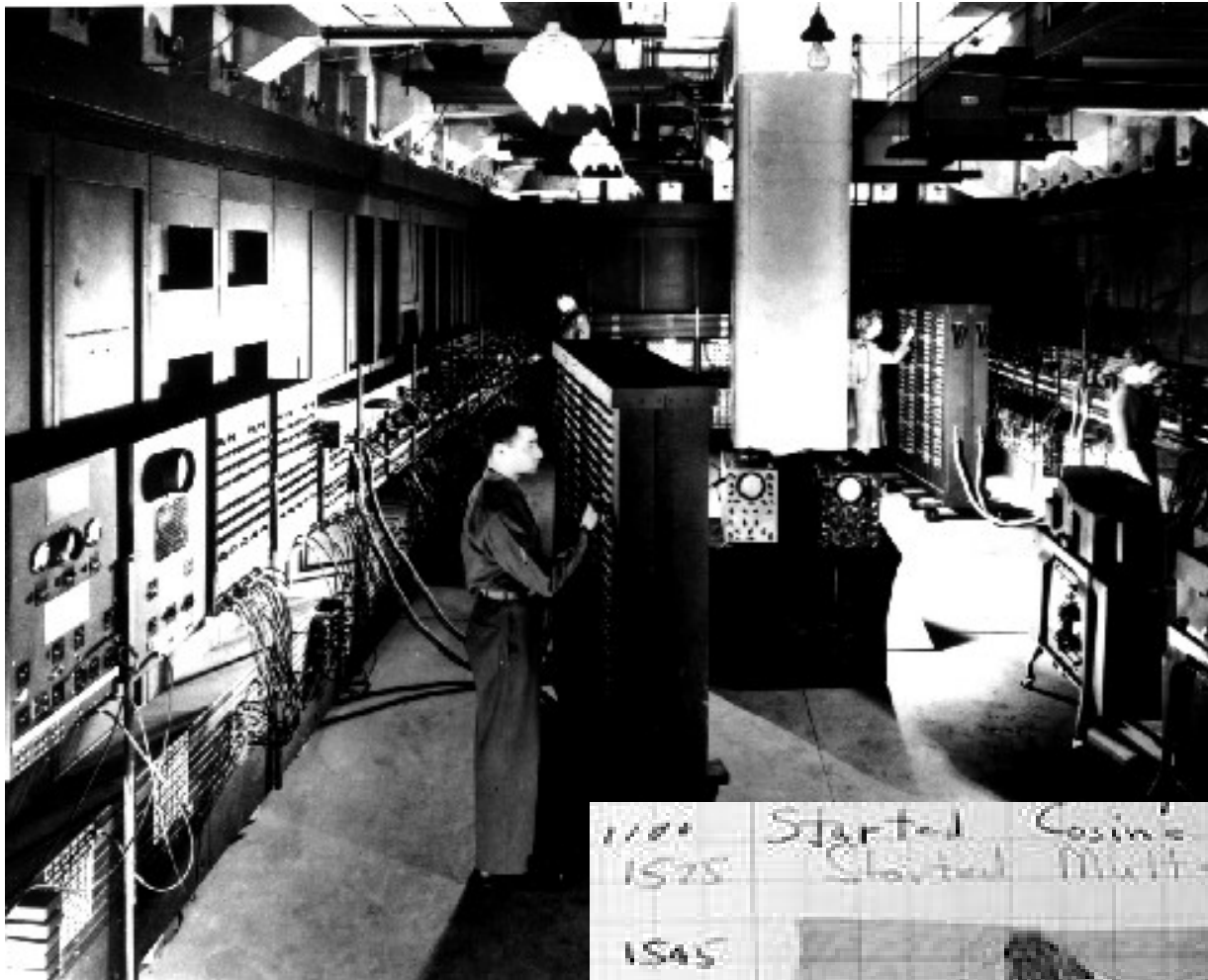
Where do cloud come from?




The momentum toward clouds

- Going distributed
 - For power, storage, profit
- Externalizing computing
 - The Pay-as-you-go model
- The ultimate dream
 - The power-grid metaphor

In the beginning...



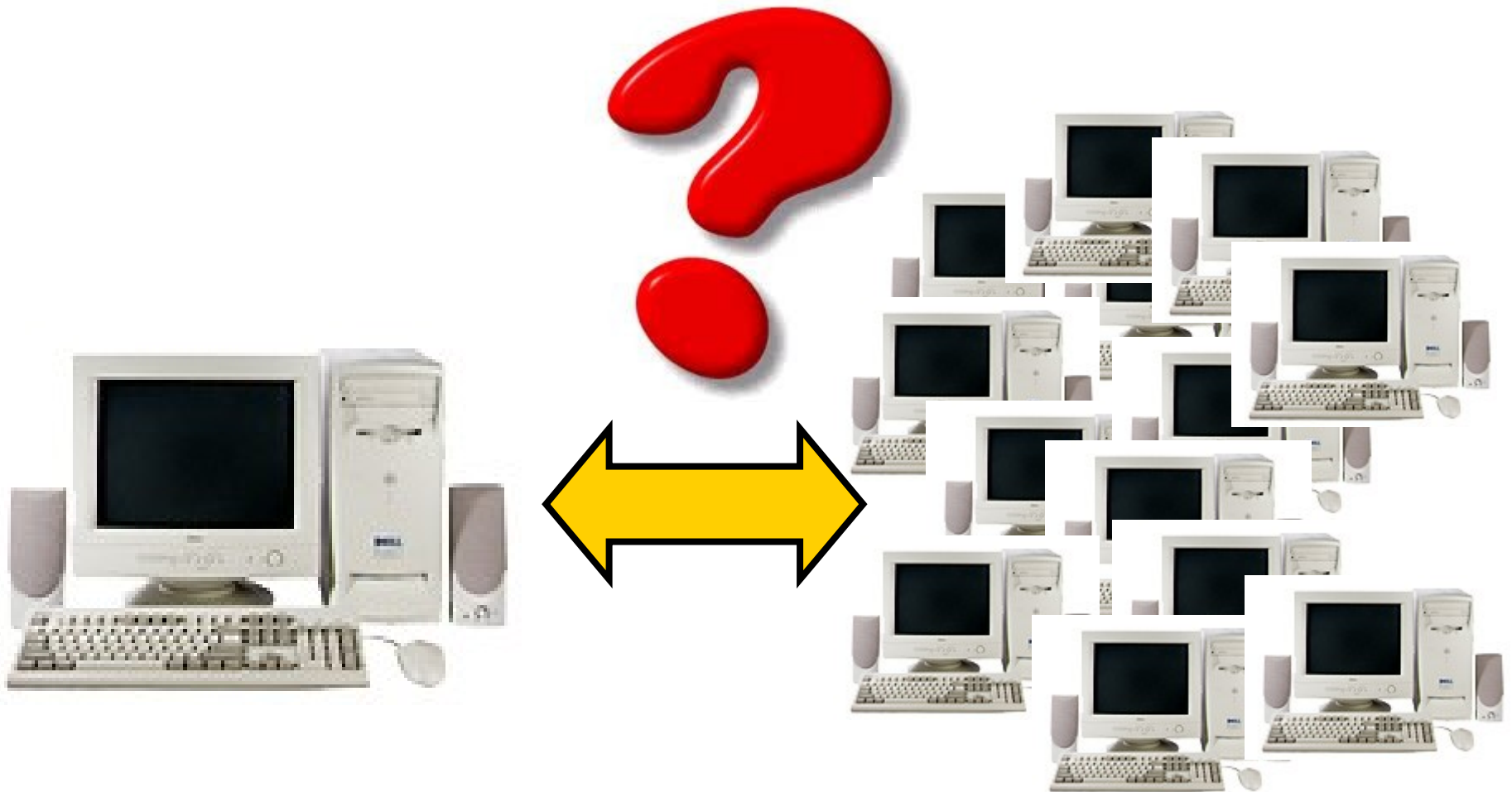
1100 Started Cosine Tape (Sine check)
1525 Started Multiplier Adder Test.
1545 Relay #70 Panel F (moth) in relay.

First actual case of bug being found.
~~1600~~ 1600 arduous started.
1700 closed down.

50 years later...



1100 Started Cosine Tape (Sine check)
1525 Started Multi Adder Test.
1545 Relay #70 Panel F (noth) in relay.
First actual case of bug being found.
1600/1600 arduino started.
1700 closed down.

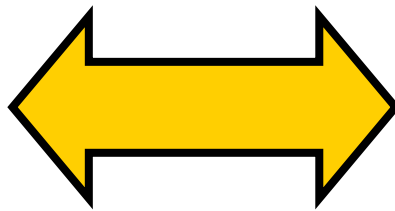
The excruciating question



Question #1: More power?

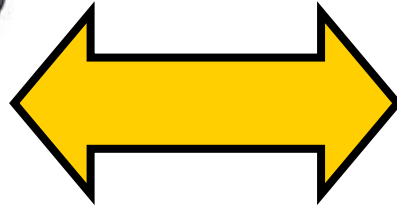


?

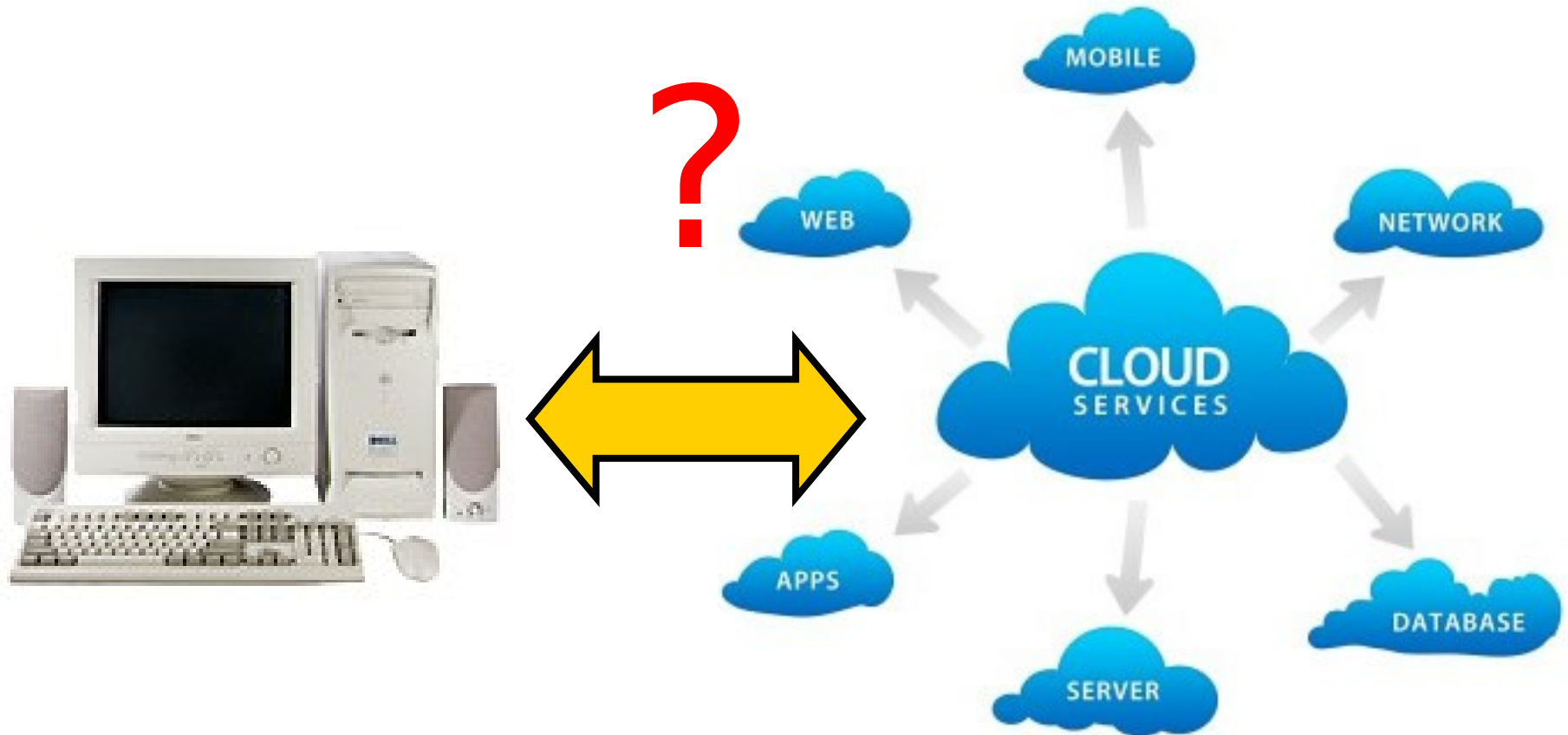


Question #2: More storage

?



Question #3: More profit



ENT de l'université de R x http://rc2009.bvdep.cor x Home - Google Docs x Gmail - All Mail - luc.bou x

https://mail.google.com/mail/?shva=1#all/p13

Google A reg IRISA Enseignement Ressources Dicos Perso EuroPar IPDPS Haltes Doc Cirque

Allow Gmail (mail.google.com) to open all email links? [Learn more](#)

+Luc Search Images Maps YouTube News Gmail Documents Calendar More

Google Luc Bougé 1 + Share

Click here to enable desktop notifications for Gmail. [Learn more](#) [Hide](#)

Gmail 1,201-1,300 of 3,721

COMPOSE

Inbox
Starred
Important
Chats
Sent Mail
Drafts
All Mail
Spam

Invite a friend +

<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Computational Science	Call For Papers: Paper Submission Deadline (March 12), Conferences in Coi	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	bnicolae	[Mapreduce-commits] r215 - in deliverables: . D4.1 D4.1/figures - Author: bn	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012	[Euro-Par 2012] Review for paper #1569567203 confirmed - Dear Prof. Andrzej	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012	[Euro-Par 2012] Review for paper #1569568283 confirmed - Dear Dr. Erik Bom	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	AlloCiné Séances	Luc, le programme ciné de la semaine - Si vous n'arrivez pas à visualiser les i	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012	[Euro-Par 2012] Review for paper #1569568309 completed - Dear Dr. Richard C	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012 (3)	[Euro-Par 2012] Review for paper #1569570517 confirmed - Dear Prof. Marco I	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Laurent Perraudeau	Réunion de préparation du conseil de laboratoire de l'Irisa - Bonjour, un cons	Feb 29
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012	[Euro-Par 2012] Review for paper #1569568111 confirmed - Dear Dr. Lee Gillar	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012	[Euro-Par 2012] Review for paper #1569564941 confirmed - Dear Mr. Bahman	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012 (3)	[Euro-Par 2012] Review for paper #1569567893 confirmed - Dear Prof. Marco I	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Patrick Senac	[asr-forum] Appel à participation ResCom 2012 - Chers Collègues, Nous vous	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Bastien Padeloup	Re: XTRA: références - Thomas a envoyé ça avec ses références à lui, vous av	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012 (2)	[Euro-Par 2012] Review for paper #1569567853 confirmed - Dear Prof. Marco I	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Rizos Sakellariou	[Euro-Par 2012] Review request for #1569567853 - Dear Dr. Wei Zheng: The pa	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012	[Euro-Par 2012] Review for paper #1569567943 declined - Dear Mr. Nikolaos K	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Euro-Par 2012 (2)	[Euro-Par 2012] Review for paper #1569562687 confirmed - Dear Prof. Marco I	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Rizos Sakellariou	[Euro-Par 2012] Review request for #1569562687 - Geia sou Giwrgo, Tha se en	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	costan	[Blobseer-publis-commits] r1733 - MapReduce2012 - Author: costan Date: 201	Feb 28
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	ACM MemberNet	Today's Topic: ACM MemberNet - Tuesday, February 28, 2012 MemberNet Wel	Feb 28

ENT de l'université de R x http://rc2009.bvdep.cor x Home - Google Docs x Convention INSA/ENS - C x

Google Agreg IRISA Enseignement Ressources Dicos Perso EuroPar IPDPS Haltes Doc Cirque

Convention INSA/ENS ☆

File Edit View Insert Format Tools Table Help

Comments Share

Luc Bougé ▾

Heading 1 Arial 24 B I U A A

1 2 3 4 5 6 7

Le courrier suivant résume les échanges de services entre l'INSA Rennes et l'ENS Cachan/Rennes en ce qui concerne l'informatique et formule des propositions pour gérer ces échanges dans un cadre contractuel bien défini.

État des lieux

Arnaud Jobin

Arnaud Jobin a obtenu un monitorat de l'INSA Rennes à la rentrée 2008. Il a effectué des enseignements dans les cursus de l'ENS en 2009-2010 et 2010-2011.

- En 2009-2010, 24h eqTD dans le cours d'algorithmique de 1er semestre de 1e année + 6h eqTD en compléments de mathématiques, soit 30h eqTD.
- En 2010-2011, 24h eqTD dans le cours d'algorithmique de 1er semestre de 1e année.

Les heures 2009-2010 ont été remboursées à l'ENS au tarif des heures complémentaires (environ 40 Euros/h) par une convention signée à l'été 2010. Par contre, les heures 2010-2011 restent non soldées.

Amélie Stainer

Amélie Stainer, en contrat doctoral avec l'Université Rennes 1 depuis la rentrée 2010, a obtenu une mission d'enseignement de l'INSA Rennes à la rentrée 2010. Dans ce cadre, elle

The key idea: externalization

- Motherhood idea: Buy the service, not the infrastructure necessary to produce it
 - Machines
 - People
 - Experience
- Pushed by IBM, late '90s
 - Various interpretations!
- Then, reincarnate as grid computing
- Currently, cloud computing
- And tomorrow?



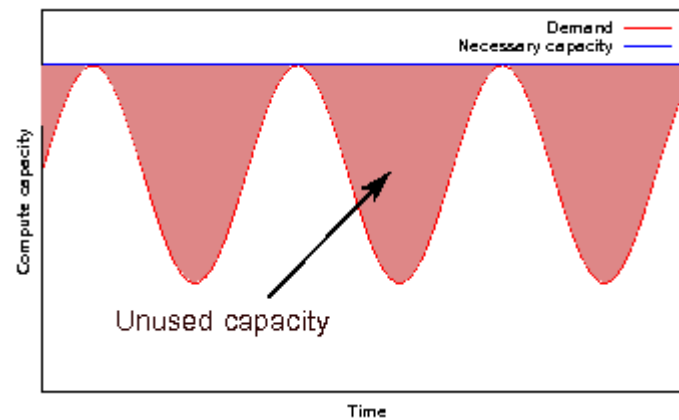
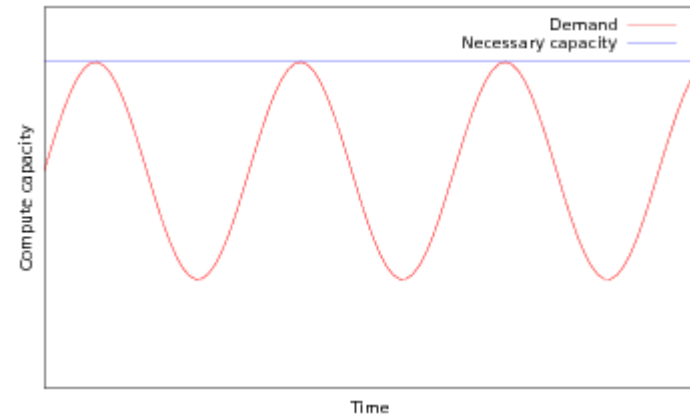
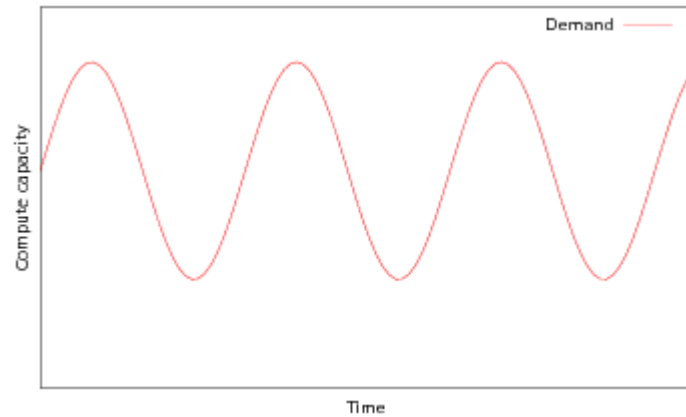
On Demand Computing

Technologies
and Strategies

Craig Fellenstein



The basic picture



Externalizing to the cloud

- Keep some small set of resources in-house
 - Safety
 - Competence
- Request resources from the cloud
 - On-demand, real-time
 - Pay-as-you-go pricing model
 - Do not support any fixed cost
- Service-Level Agreement Guaranteed by contract
 - Various level of offers
- No long-term commitment to any provider
 - Regular economic laws apply

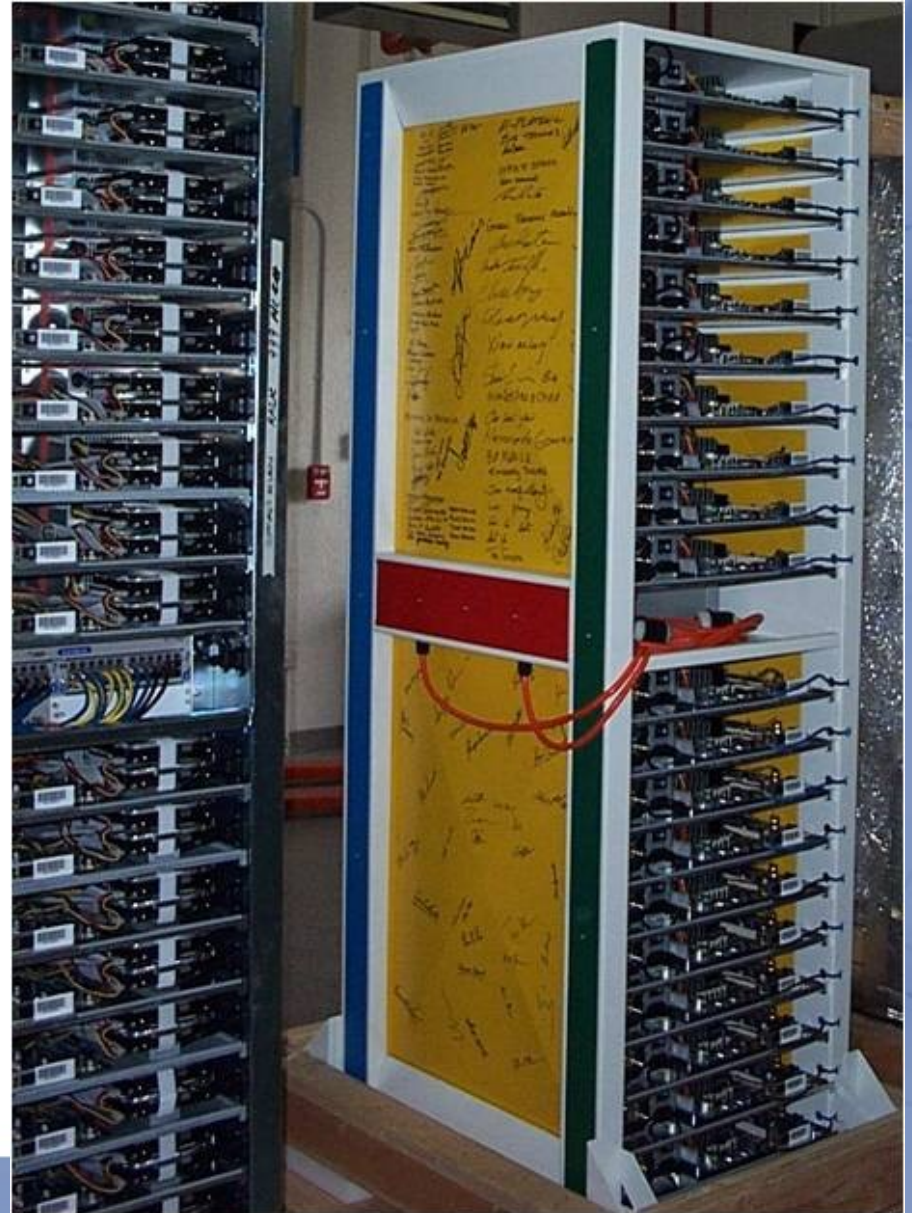


Google cluster, 1997

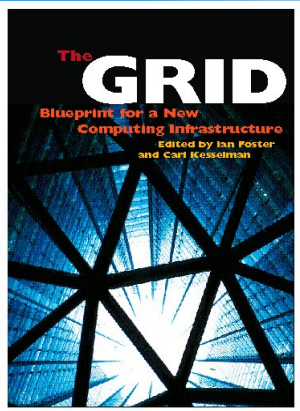
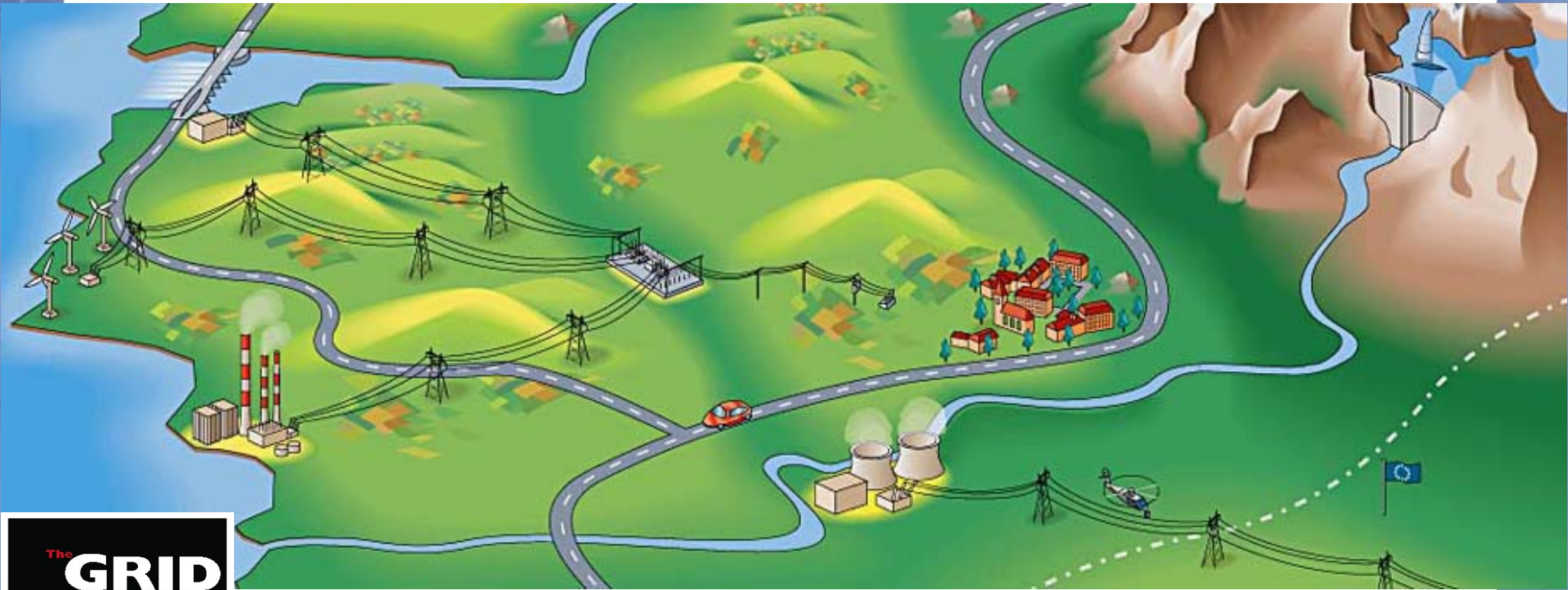


Google cluster, (almost) today

- 36 data centers
 - > 800K servers
- 40 servers/rack

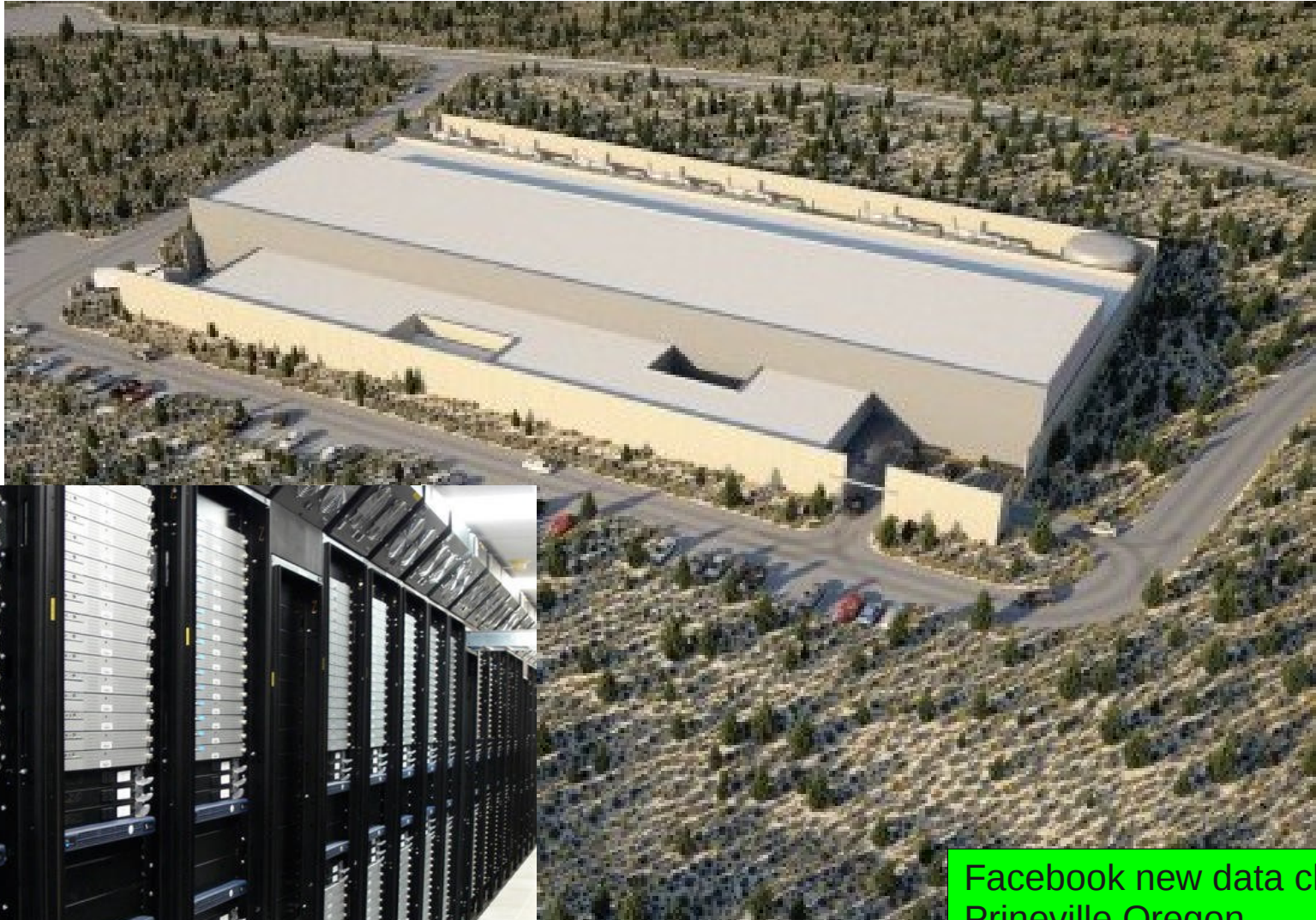


The Power-Grid Metaphor



Ian Foster, Carl Kesselman, ca. 1995

Cloud computing technology



Facebook new data cluster,
Prineville Oregon.
30 PB storage

What made cloud possible?

- Clusters
 - Opening the way toward distributed computing for non-distributed task
- Grids
 - Large-scale, heterogeneous computing
- Virtual machines
 - Hiding the hardware altogether
- High-speed networks
 - Hiding the location altogether
- Hardware packaging and power management
- OK, but why not earlier?

Key #1: Clusters



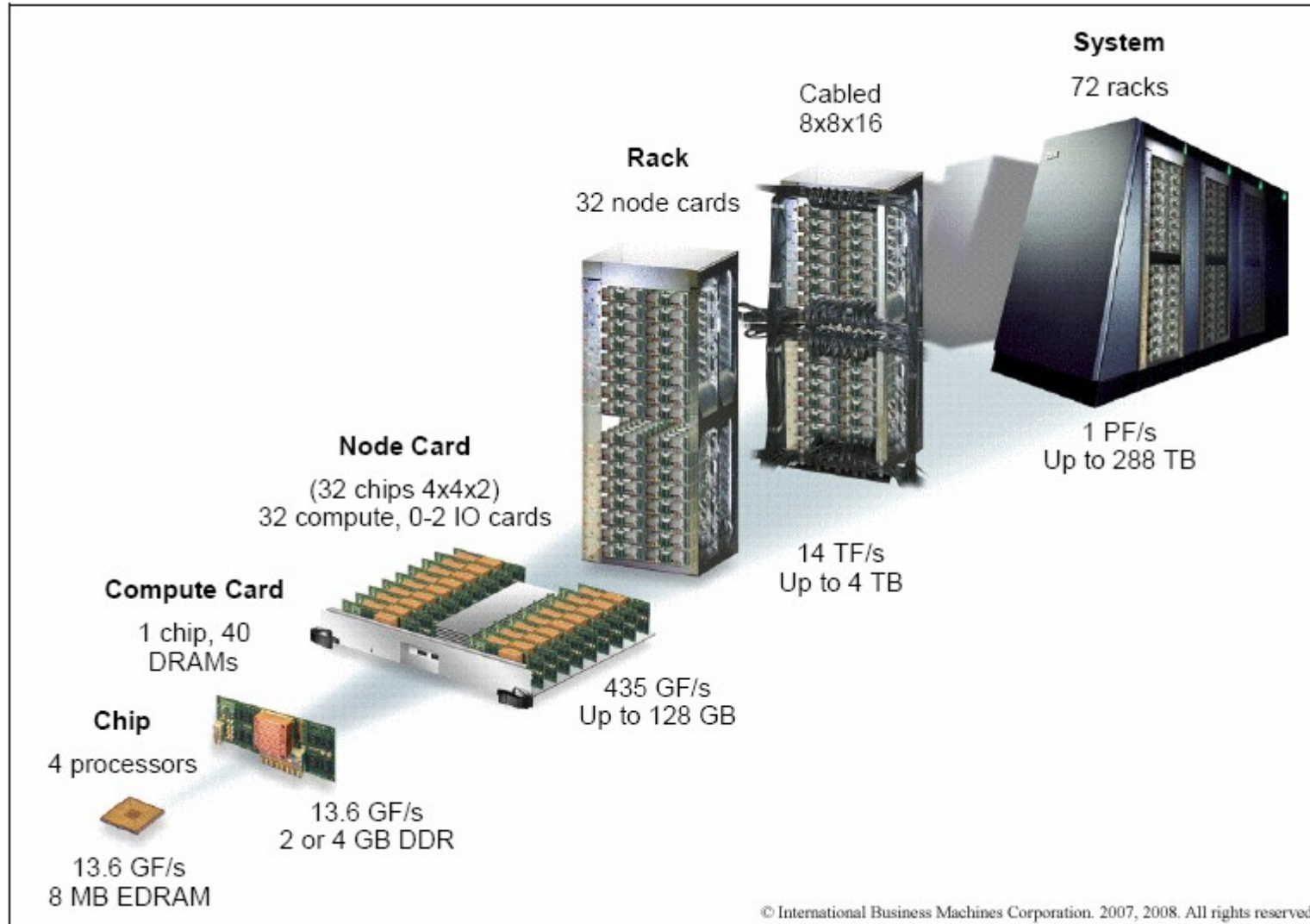
Beowulf cluster, Thomas Sterling, 1994

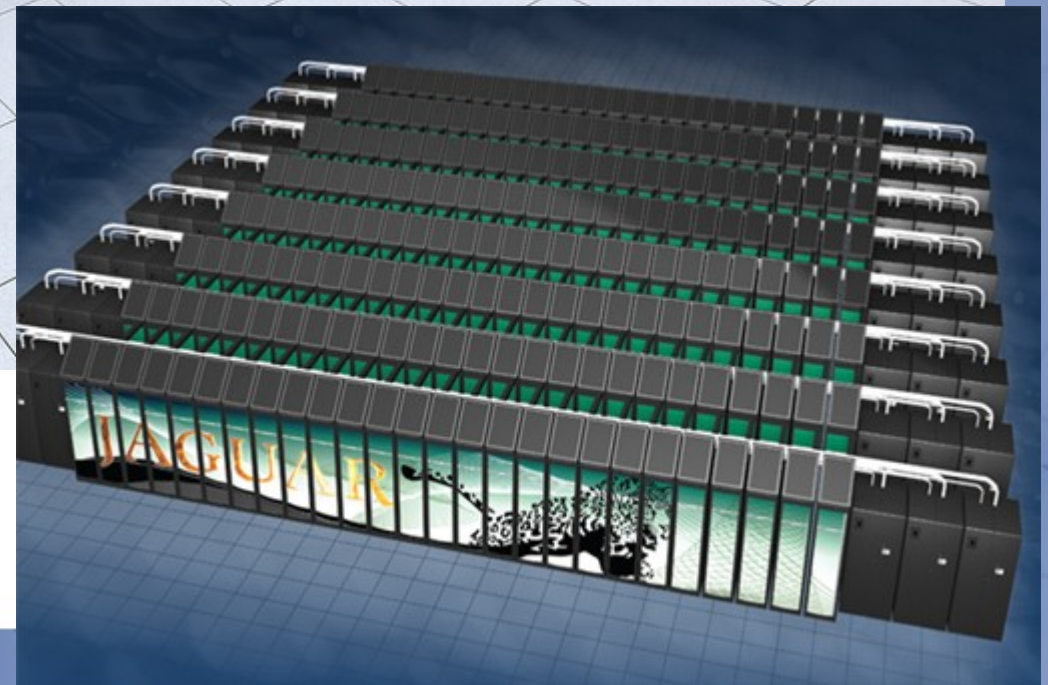


Clusters today, in France



Very large clusters in the world

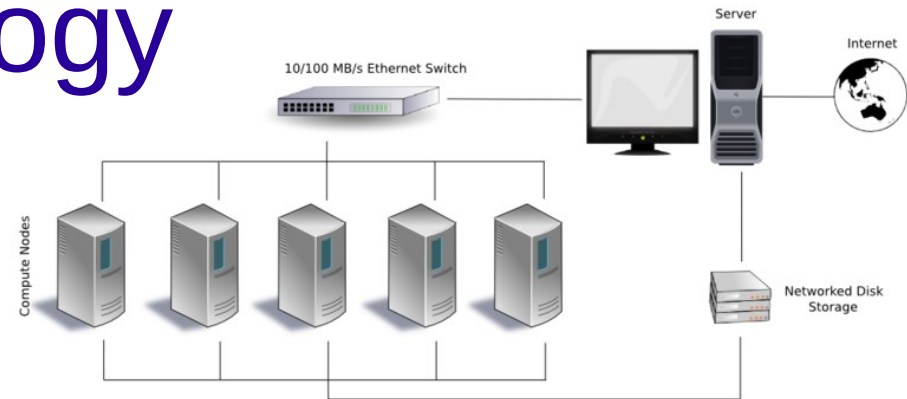




Jaguar supercomputer, TOP500 #3,
224,000 processing cores,
each with 2 GB of local memory.

Clusters technology

- Data sharing
 - Distributed file system
- Message passing and communication
 - MPI
- Task scheduling
 - Node failure management
 - Integrated failure recovery mechanism
- Debugging and monitoring
- Operating system
 - Linux, Microsoft
 - SSI approach: Mosix, Kerrighed



Applications

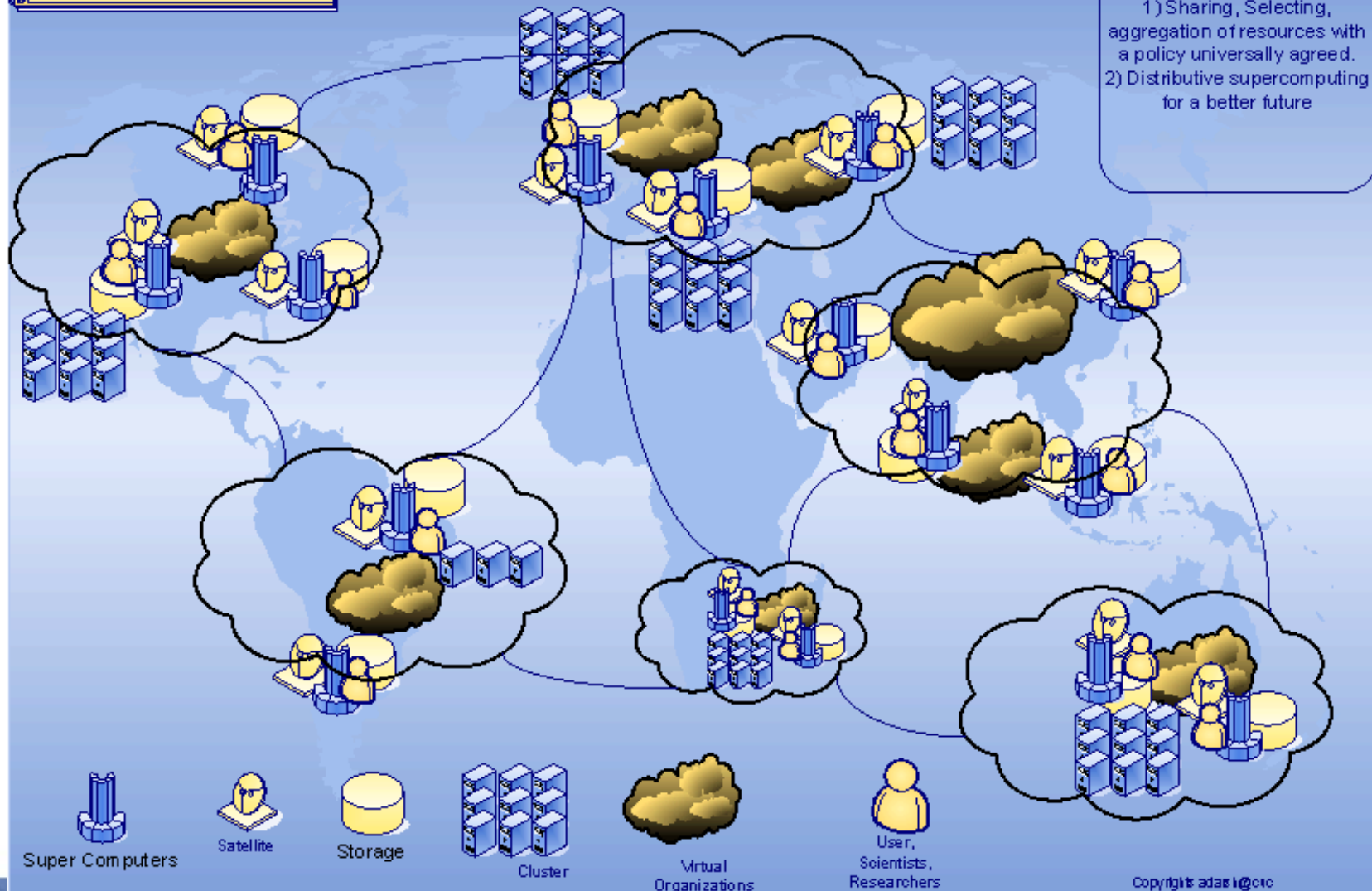
- Departmental clusters
 - Specific hardware
 - Poor-man's supercomputer: cycle-stealing
- Many traditional applications
 - Data bases: Oracle
 - Numerical crunching
 - Imaging: virtual reality



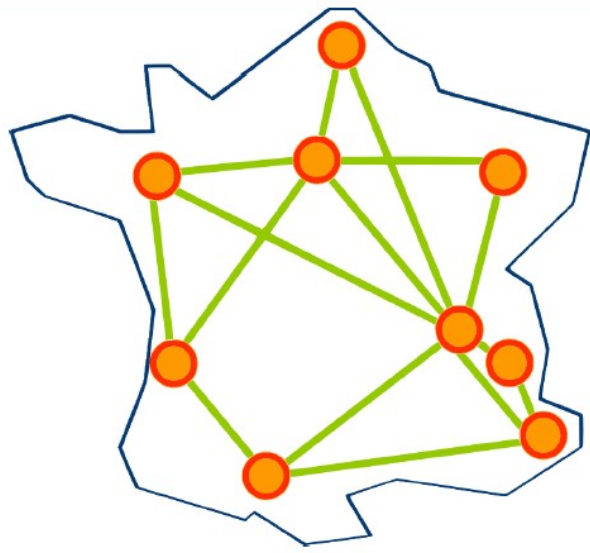
Key #2: Grids

A federation of clusters

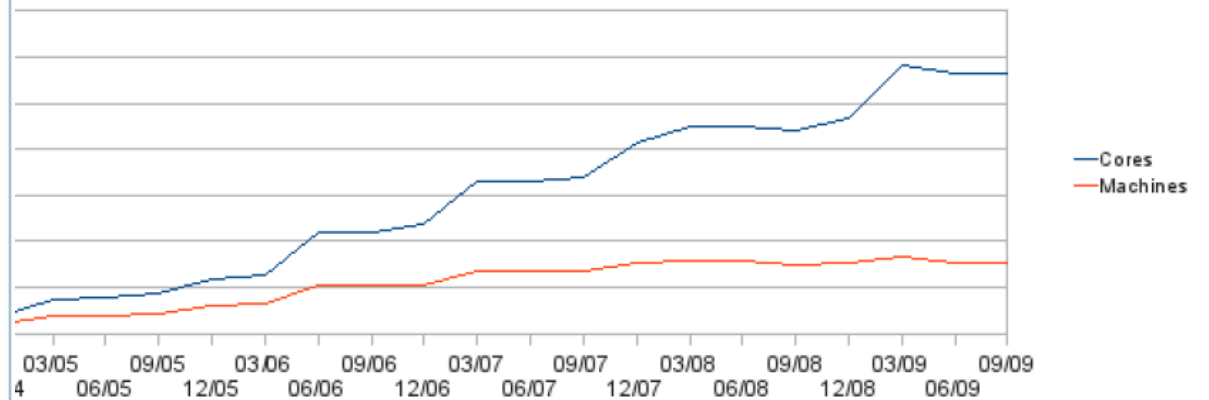
Grid Computing



Grid5000



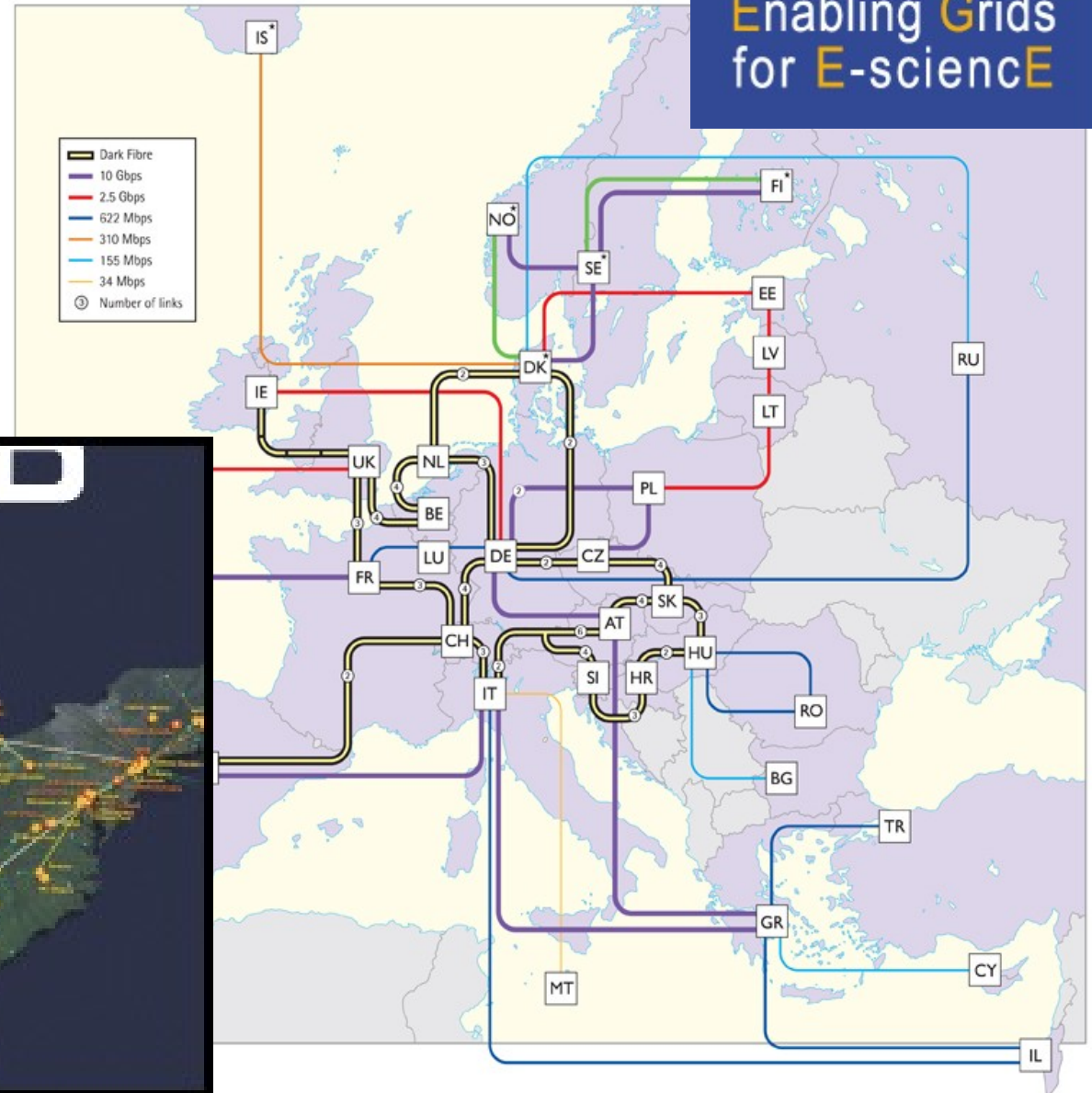
Évolution matérielle de Grid'5000



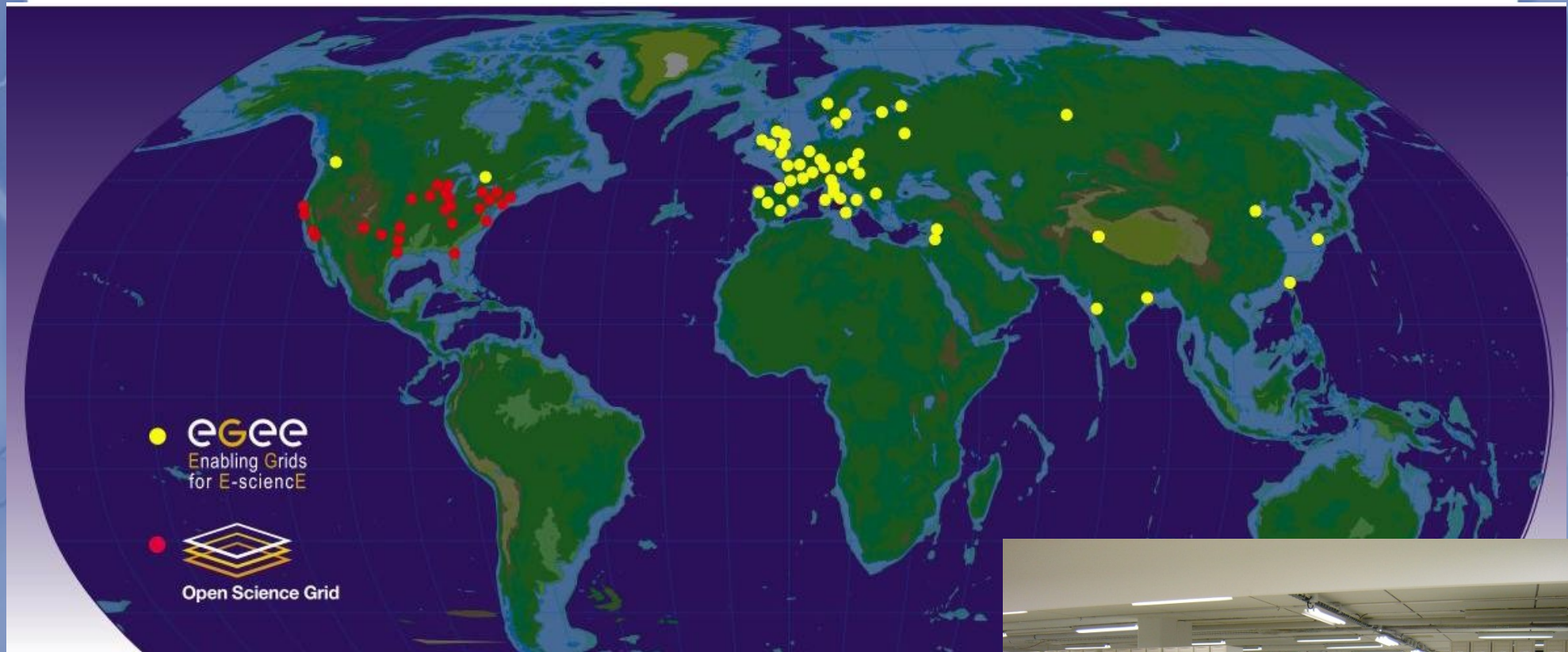
And elsewhere in the world



eGEE
Enabling Grids
for E-science



Grid at CERN



Grid technology

- Grids = clusters
 - Size, heterogeneity, load
- Additional specificity
 - Computing resources are not administered centrally
 - Open standards are used
 - Nontrivial quality of service is achieved
- Virtual organization
 - Sharing power, data, but also resources and people
- Key problems:
 - External user interface: a single virtual organization
 - Security in spite of multiple organizations
 - Failure resiliency



Key #3: Virtual machines

- Completely isolated guest operating system installation within a normal host operating system
 - Software emulation
 - Hardware virtualization
 - (in most cases) both together
- Two early examples
 - Java Virtual Machine
 - Grid 5000 approach: rebootable nodes
- Microprocessor progress: virtualization in hardware
 - Reasonably efficient execution speed

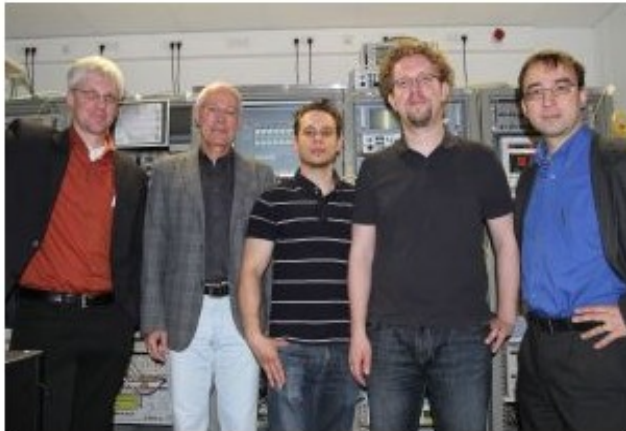


Key #4: High-speed networks

Press Release 084/2011

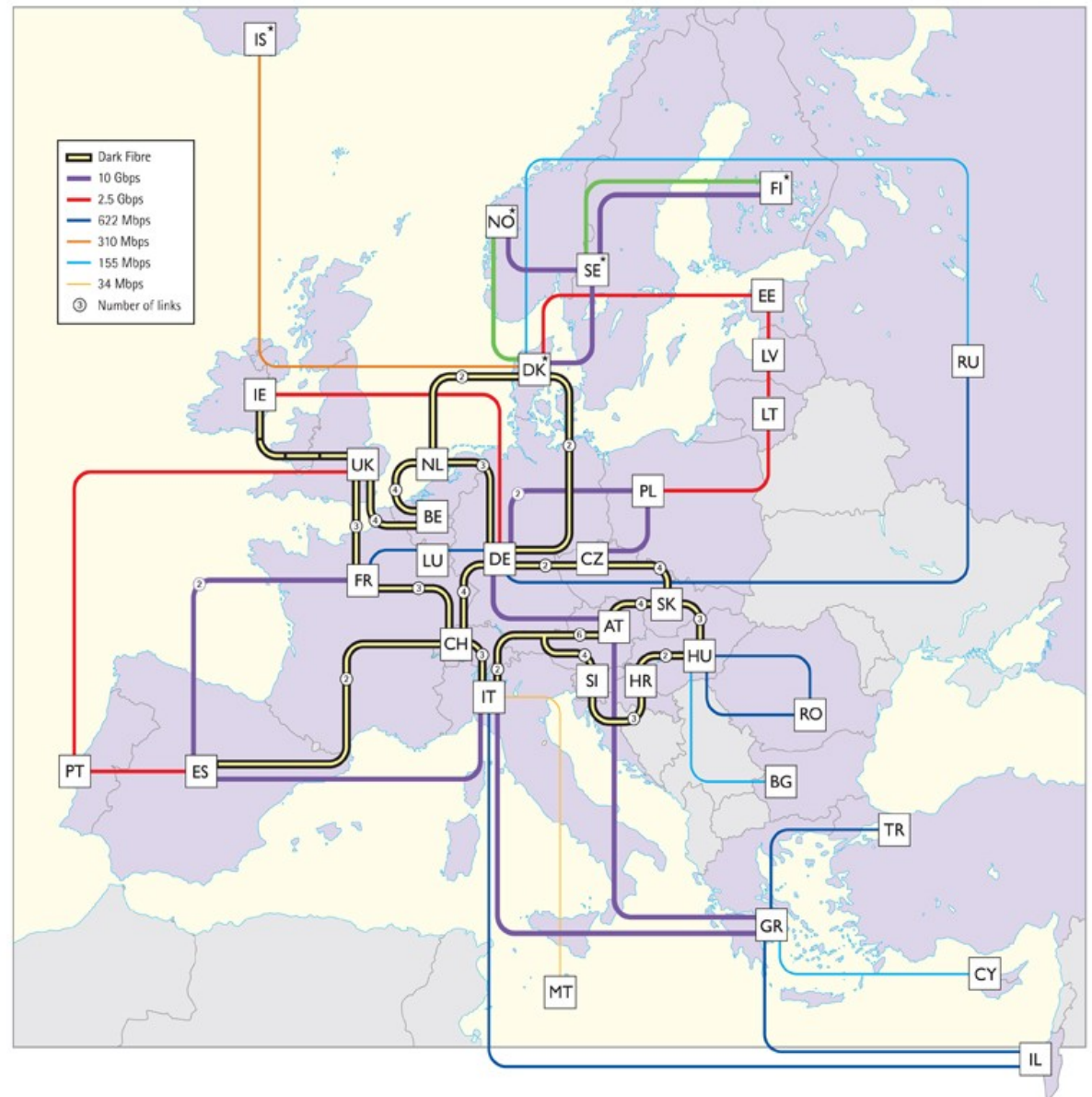
World Record in Ultra-Rapid Data Transmission

Transfer of 700 DVDs in One Second Only – Highest Bit Rate on a La



The team of Professor Leuthold (right): David Hillerkuß, René Schmogrow, and Christian Koos (from right to left). (Photo: Gabi Zachmann)

Scientists of Karlsruhe Institute of Technology (KIT) have succeeded 26 terabits per second on a single laser beam, transmitting them over decoding them successfully. This is the largest data volume ever transmitted. The process developed by KIT allows to transmit the contents of 700 DVD renowned journal "Nature Photonics" reports about this success in its issue of 10.1038/NPHOTON.2011.74).



Key #5: Hardware packaging and power management



The race for data-centers

ARCHITECTURE



Green Mountain Data Center is Buried Underground and Cooled By Norway's Fjords

by **Bridgette Meinhold**, 12/26/11

filed under: *Architecture, Sustainable Building*

Like 10

INHABITAT PHOTO GALLERY

NEXT IMAGE ►



Buried deep underneath the mountains near Stavanger, Norway, the **Green Mountain Data Center** is quite possibly the greenest data center in the world. Powered by renewable energy from nearby sources and cooled with water from the adjacent fjord, the cavernous data center is all about energy efficiency. The center dramatically reduces its cooling costs and energy use by tapping into the 8 degree Celsius water from the fjord.

Seamlessly shifting workloads between data centers might lead to the management possibility being called “follow the moon” which takes advantage of lower costs for power and cooling during overnight hours. Virtualized workloads would be shifted across data centers in different time zones to capture savings from off-peak utility rates. Go man, go!



OK, but why not earlier?

- Companies/organizations were not ready to entrust data to a foreign organization
 - Confidentiality
 - Legal problems with respect to actual data storage location
 - Slow transmissions
 - Lack of tools for fine monitoring
- Intensive ad campaigns from cloud companies
 - Amazon: economic arguments
 - Google and Yahoo! example
- Feeling that entrusting is unavoidable
 - “Data deluge”



Cloud computing





enomalism Elastic Computing Platform

2.0

Amazon Web Services @ Amazon.com

Home /

BackForward Reload Stop Home
Latest Headlines CASudRA
Google large scale da
Amazon Web Services @ A

Coming

amazon.com

- Be
- Cli

Shop All Departments

Make Money

Enomalis
Elastic C
flexible
without



Key Ben

- AP
- Sc
- Ap
- Be
- Su
- Su
- Im
- Be
- Im
- Ac
- Se
- Bil
- Di
- Liv
- Vir

About AWS

- [Why Use AWS?](#)
- [In the News](#)
- [Upcoming Events](#)
- [Customer Case Studies](#)
- [Solutions Catalog](#)
- [Partners](#)

[Careers at AWS](#)

Infrastructure Services

- [Amazon Elastic Compute Cloud](#)
- [Amazon SimpleDB](#)
- [Amazon Simple Storage Service](#)
- [Amazon Simple Queue Service](#)

[AWS Premium Support](#)

Payments & Billing Services

Done

IBM

IBM Software Strategy Group

IBM Google Announcement on Internet-Scale Computing ("Cloud Computing Model")

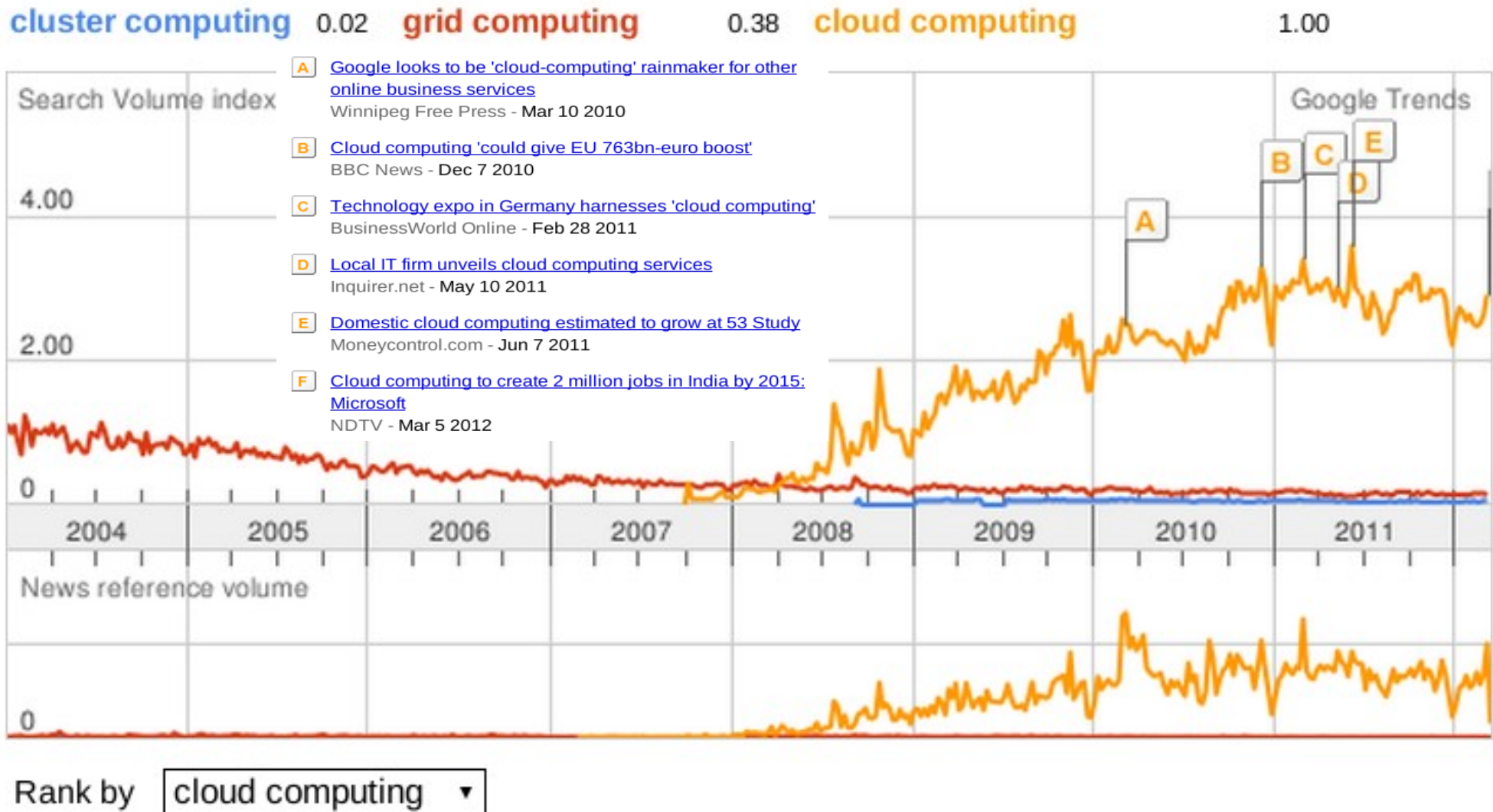
Oct 8, 2007

HiPODS

IBM Confidential until Oct 8, 2007

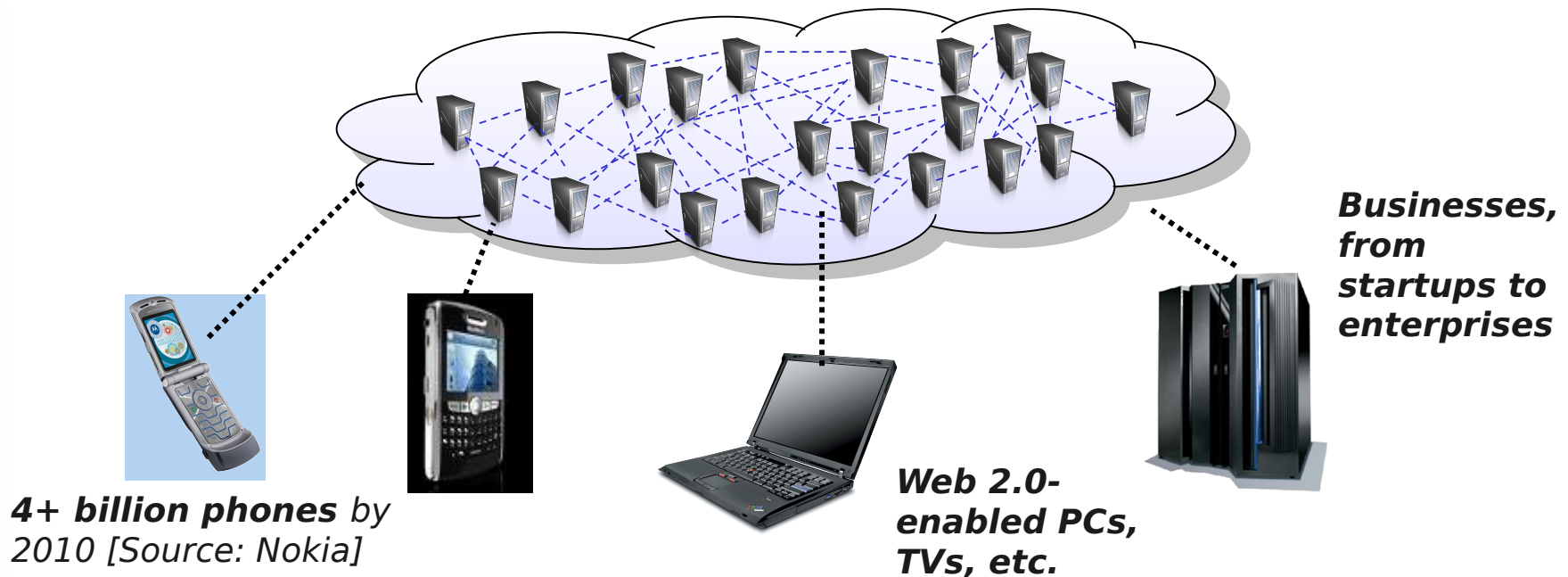
© 2007 IBM Corporation

Beware of the Cloud Hype !



What is cloud computing?

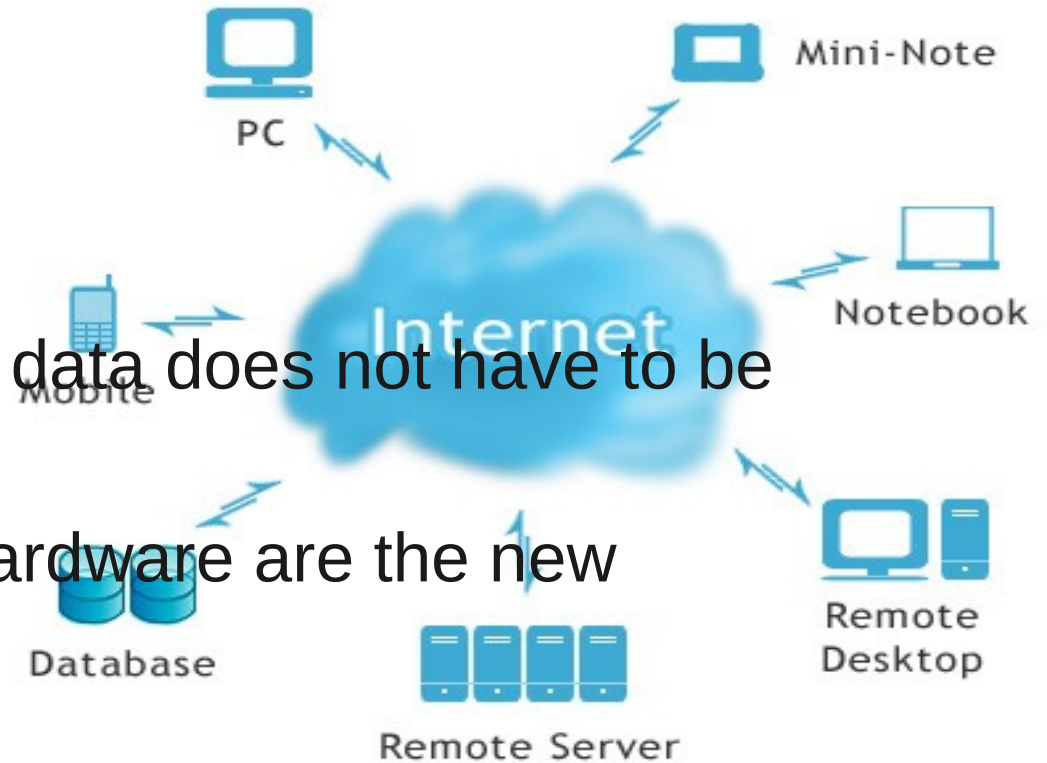
An emerging computing paradigm where data and services reside in massively scalable data centers and can be ubiquitously accessed from any connected devices over the internet.



Credit: IBM Corp.

Key concepts

- Processing 1000x more data does not have to be 1000x harder
- Cycles and bytes, not hardware are the new commodity
- Cloud computing is
 - Providing services on virtual machines allocated on top of a large physical machine pool
 - A method to address scalability and availability concerns for large scale applications
 - Democratized distributed computing



Cloud functionality

- SaaS: Software as a Service
 - Google Mail, Google Docs
- DaaS: Data as a Service
 - Cloud as a data repository
- PaaS: Platform as a Service
 - Amazon, Azure: select your VM
- IaaS: Infrastructure as a Service
 - Grid 5000: manage your own VM
- HaaS: Hardware as a Service



Online Utilities and Applications Generally Referred to as "Cloud Computing" or Software as a Service (SaaS)

Provider	Utilities			Services			Applications	
	Network Utilities	Online storage	Online processing	Developer Environment	Application Services & Tools	Business Process Outsourcing	Online Enterprise Applications	Online Consumer Applications
Akamai	Web Acceleration							
Amazon		S3	EC2		SimpleDB & SQS			
AT&T	√	√	√		AT&T Web Meeting			
Box.net		√			√			√
Google		Google Apps Engine, Android, Open Social					Google Apps Gmail/Postini	Google Apps Gmail
IBM	√	Blue Cloud						
Microsoft		√	√	Azure			Office Online	Office Online Hotmail
Salesforce.com				Force.com Links to Google Apps and Amazon Web Services		√	√	
SAP					SAP Web Application Library	√	√	
Sun Microsystems	√	Metro Web Services		Sun SOA Java Dev.Pack			Open Office	
Terremark Worldwide	√	√	√					
Yahoo!				Pipes				Yahoo! Mail
Xcalibre	√	√	√					

A zoom on Amazon

A set of APIs and business models which give developer-level access to Amazon's infrastructure and content:

📦 Data As A Service

- 📦 Amazon E-Commerce Service
- 📦 Amazon Historical Pricing

📦 Search As A Service

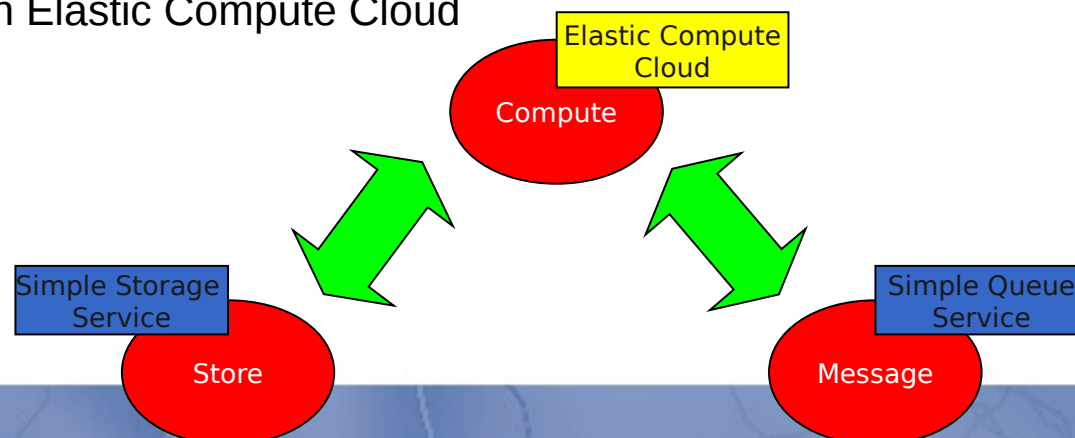
- 📦 Alexa Web Information Service
- 📦 Alexa Top Sites
- 📦 Alexa Site Thumbnail
- 📦 Alexa Web Search Platform

📦 Infrastructure As A Service

- 📦 Amazon Simple Queue Service
- 📦 Amazon Simple Storage Service
- 📦 Amazon Elastic Compute Cloud

📦 People As A Service

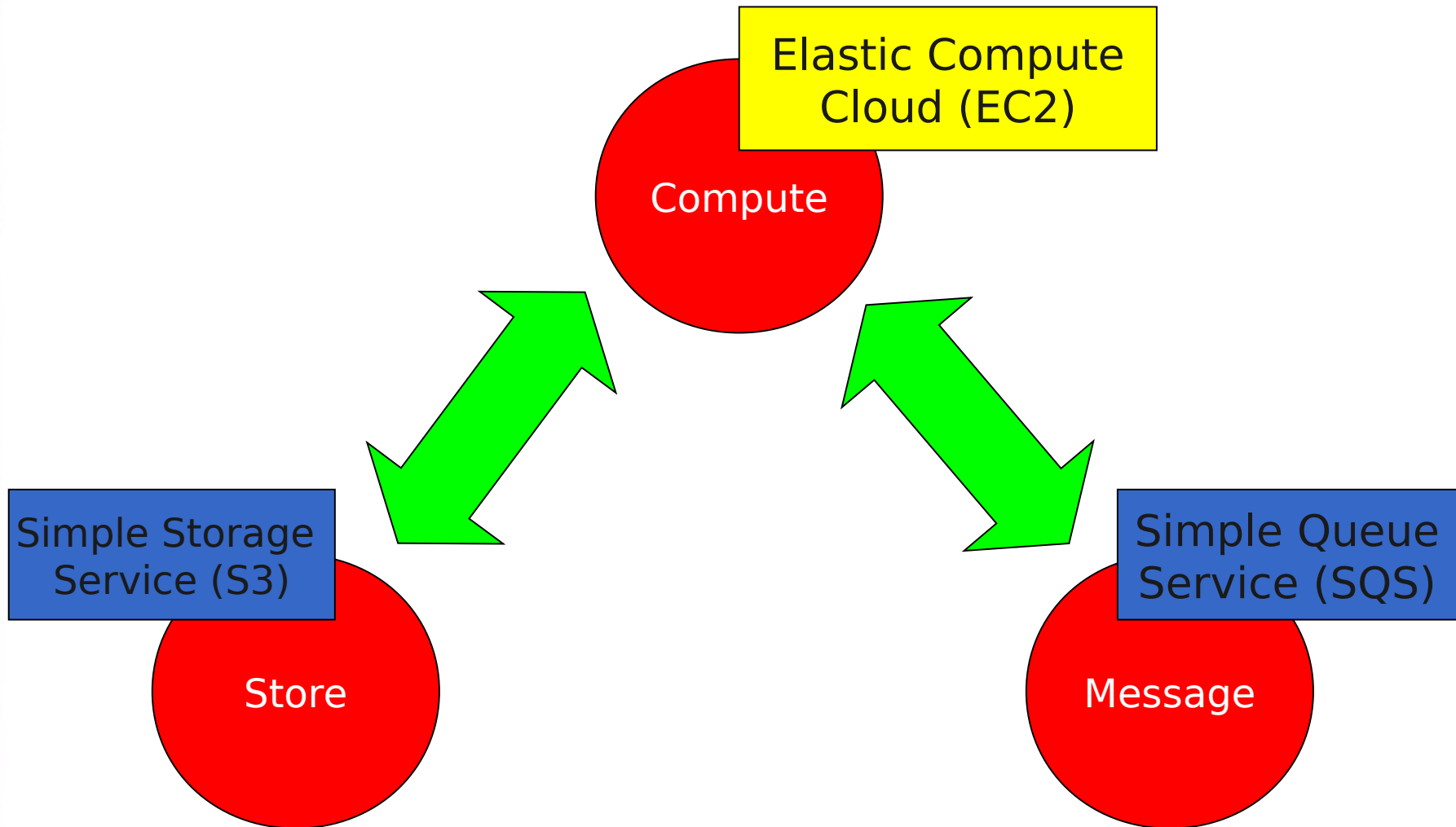
- 📦 Amazon Mechanical Turk



Credits: Jeff Barr, Amazon



Amazon Web services



Region:	EU (Ireland) ▼	
	Linux/UNIX Usage	Windows Usage
Standard On-Demand Instances		
Small (Default)	\$0.090 per Hour	\$0.115 per Hour
Medium	\$0.180 per Hour	\$0.230 per Hour
Large	\$0.360 per Hour	\$0.460 per Hour
Extra Large	\$0.720 per Hour	\$0.920 per Hour
Micro On-Demand Instances		
Micro	\$0.025 per Hour	\$0.035 per Hour
Hi-Memory On-Demand Instances		
Extra Large	\$0.506 per Hour	\$0.570 per Hour
Double Extra Large	\$1.012 per Hour	\$1.140 per Hour
Quadruple Extra Large	\$2.024 per Hour	\$2.280 per Hour
Hi-CPU On-Demand Instances		
Medium	\$0.186 per Hour	\$0.285 per Hour
Extra Large	\$0.744 per Hour	\$1.140 per Hour
Cluster Compute Instances		
Quadruple Extra Large	N/A*	N/A*
Cluster GPU Instances		
Quadruple Extra Large	N/A*	N/A*
* Cluster Compute and Cluster GPU Instances are currently only available in the US East (Virginia) Region.		

Surcouf, March 2012:
3 TB network disk for
350€ = \$450

Amazon, March 2012:
3 TB for 3 years =
 $0.093 \times 3000 \times 36 =$
\$10,000

Storage Pricing

Region: EU (Ireland) ▼		
	Standard Storage	Reduced Redundancy Storage
First 1 TB / month	\$0.125 per GB	\$0.093 per GB
Next 49 TB / month	\$0.110 per GB	\$0.083 per GB
Next 450 TB / month	\$0.095 per GB	\$0.073 per GB
Next 500 TB / month	\$0.090 per GB	\$0.063 per GB
Next 4000 TB / month	\$0.080 per GB	\$0.053 per GB
Over 5000 TB / month	\$0.055 per GB	\$0.037 per GB

Request Pricing

Region: EU (Ireland) ▼	
	Pricing
PUT, COPY, POST, or LIST Requests	\$0.01 per 1,000 requests
GET and all other Requests †	\$0.01 per 10,000 requests
† No charge for delete requests	

Data Transfer Pricing

Region: EU (Ireland) ▼	
Data Transfer IN	
All data transfer in	\$0.000 per GB
Data Transfer OUT	
First 1 GB / month	\$0.000 per GB
Up to 10 TB / month	\$0.120 per GB
Next 40 TB / month	\$0.090 per GB
Next 100 TB / month	\$0.070 per GB
Next 350 TB / month	\$0.050 per GB
Next 524 TB / month	Contact Us
Next 4 PB / month	Contact Us
Greater than 5 PB / month	Contact Us

RSS: Designed to provide 99.99% durability and 99.99% availability of objects over a given year. This durability level corresponds to an average annual expected loss of 0.01% of objects.

Clouds: scientific challenges



Clouds: scientific challenges

- Almost no new scientific subject here...
- But one single parameter makes all subjects completely different...



Scale



3 scientific challenges

(among many...)

- Data management at a very large scale
- New programming models for very large-scale programming
- New complexity models for very large-scale computing

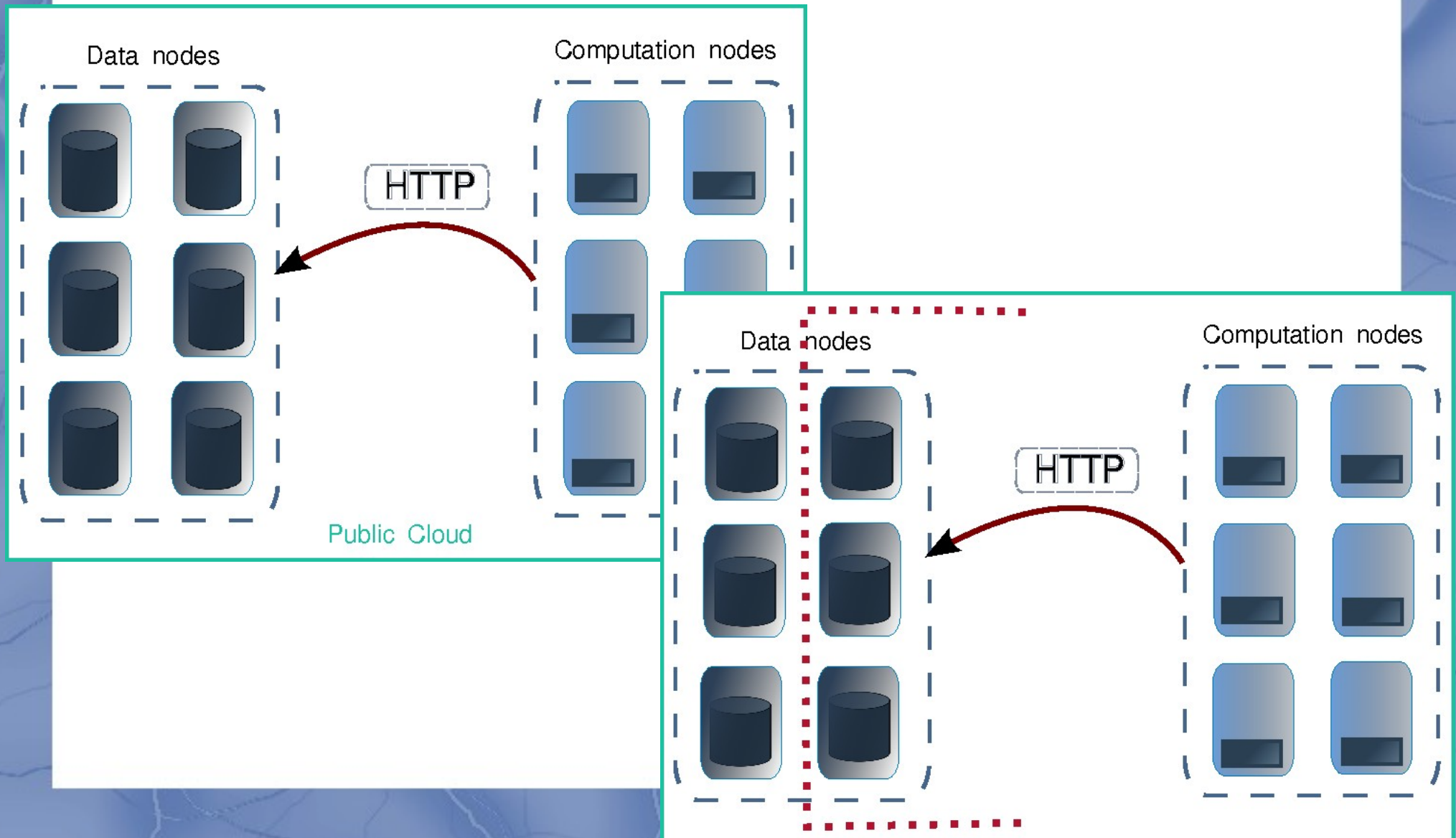


#1: Data management at a very large scale

Cloud data storage services

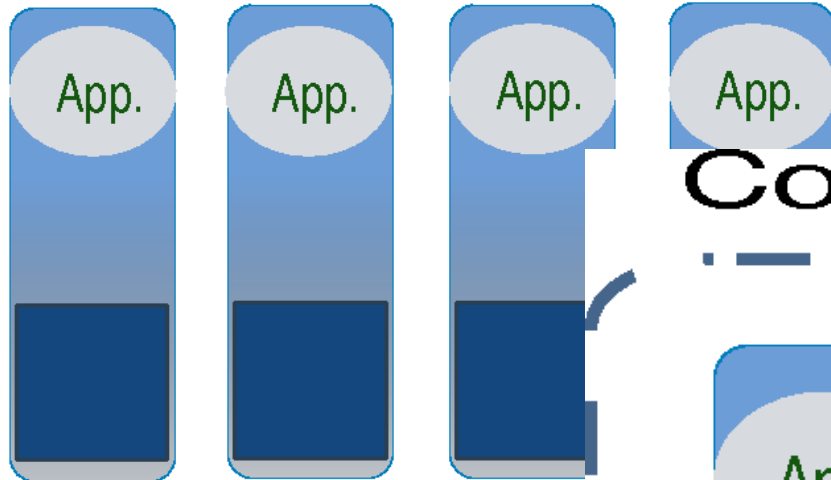
- Advantages
 - High data availability
 - Versioning
- Limitations
 - No support for concurrent accesses
 - No fine-grain data access
 - Limited object size
 - Low throughput

The cloud vision of data management

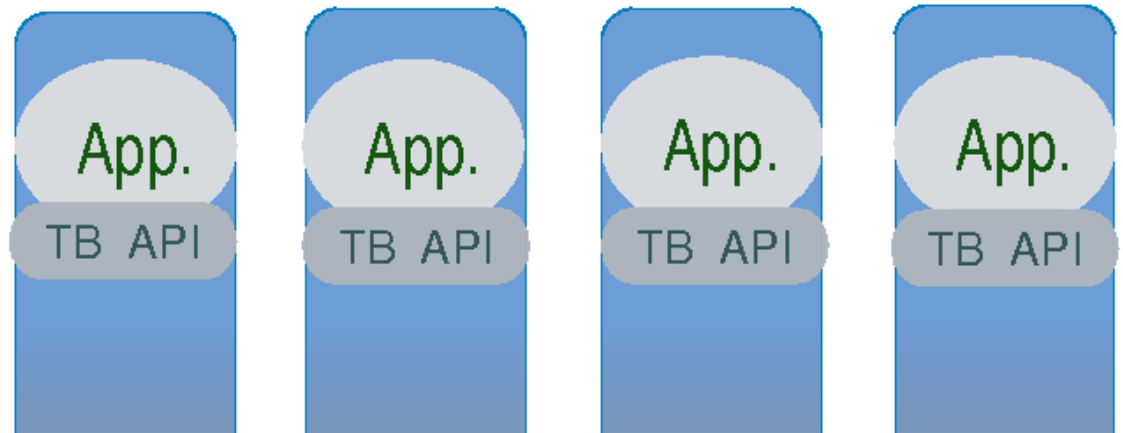


An alternative vision

Computation nodes



Computation nodes



BlobSeer

The BlobSeer approach

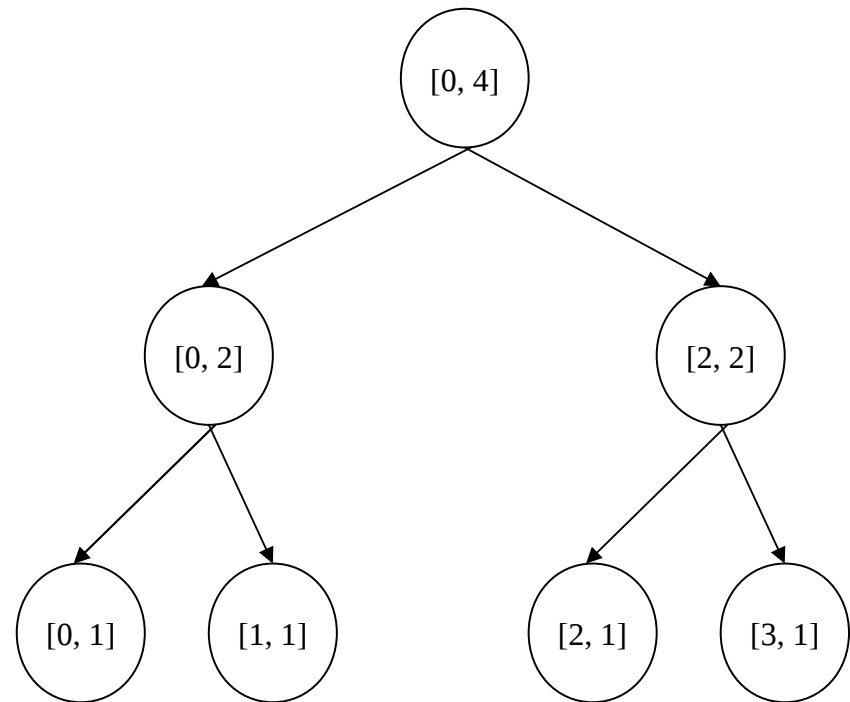
- BlobSeer: software platform for scalable, distributed BLOB management
 - Huge data (TB) - BLOBs: Binary Large Objects
 - Highly concurrent, fine-grain access (MB):
Read/Write/Append
 - Developed by the KerData Team at INRIA Rennes
- A back-end for higher-level, sophisticated data management systems
- Short term: highly scalable distributed file systems
- Middle term: storage for cloud services
- Long term: extremely large distributed databases

Scientific contribution: lock-free access

- Versioning-based concurrency control
- Update/append: generate new chunks rather than overwrite
- Metadata is extended to incorporate the update
- Both the old and the new version of the BLOB are accessible
- Lock-free approach: write-once, read-many concurrent access

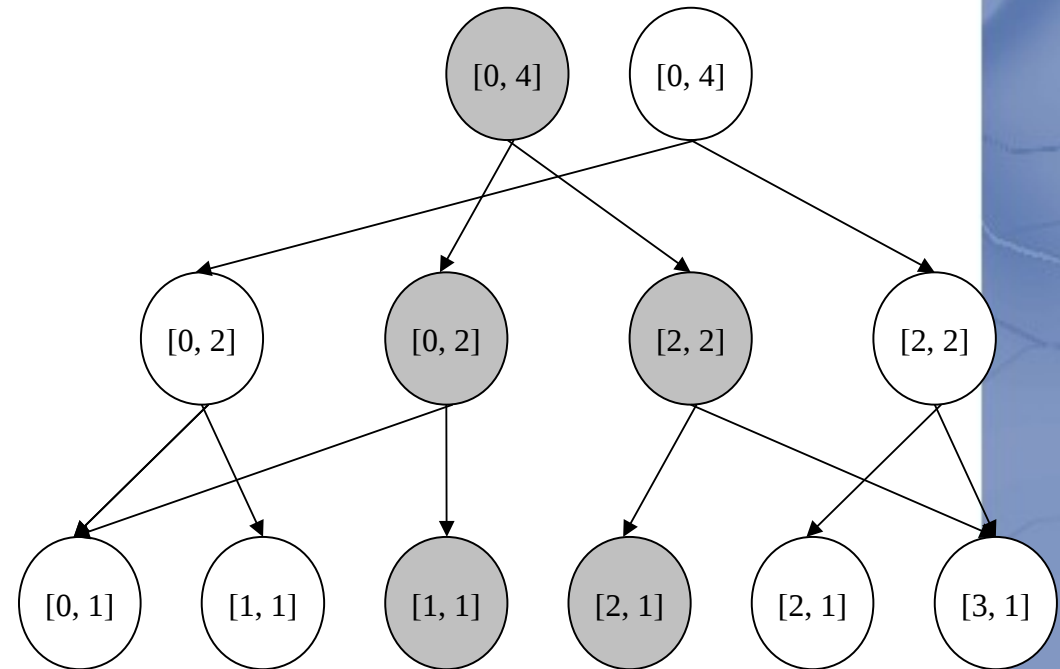
Zooming on metadata (1)

- Organized as a segment tree
- Each node covers a range of the blob identified by **[offset, size]**
- The first/second half of the range is covered by the left/right child
- Each leaf corresponds to a chunk and holds information about its location



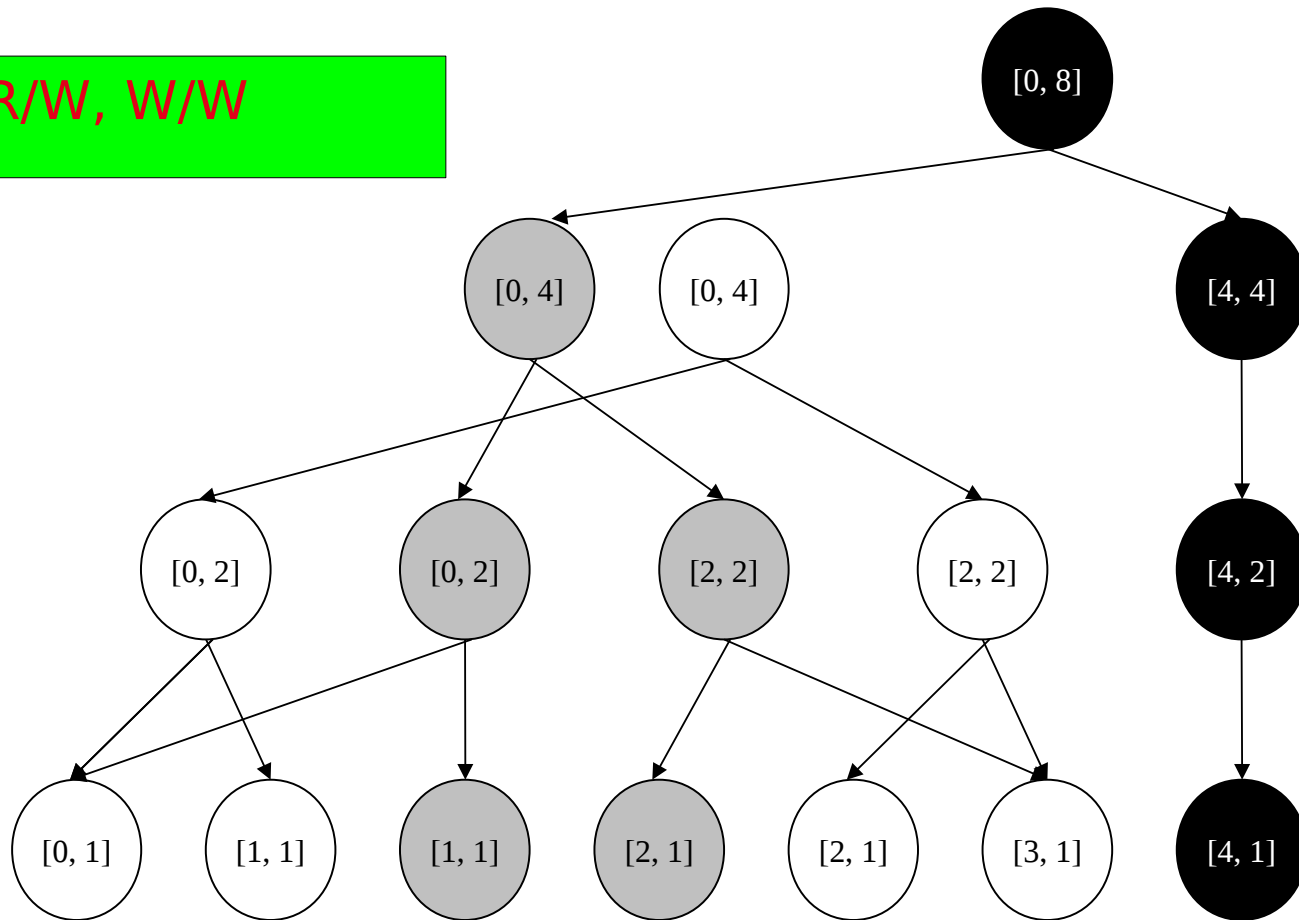
Zooming on metadata (2)

- Each node holds versioning information
- Write/Append
 - Add leaves and build subtree up to the root
 - The tree may grow one level
- Read
 - Descend from the root towards the leaves
- Tree nodes are distributed among metadata providers
- Highly scalable access concurrency:



Zooming on metadata (3)

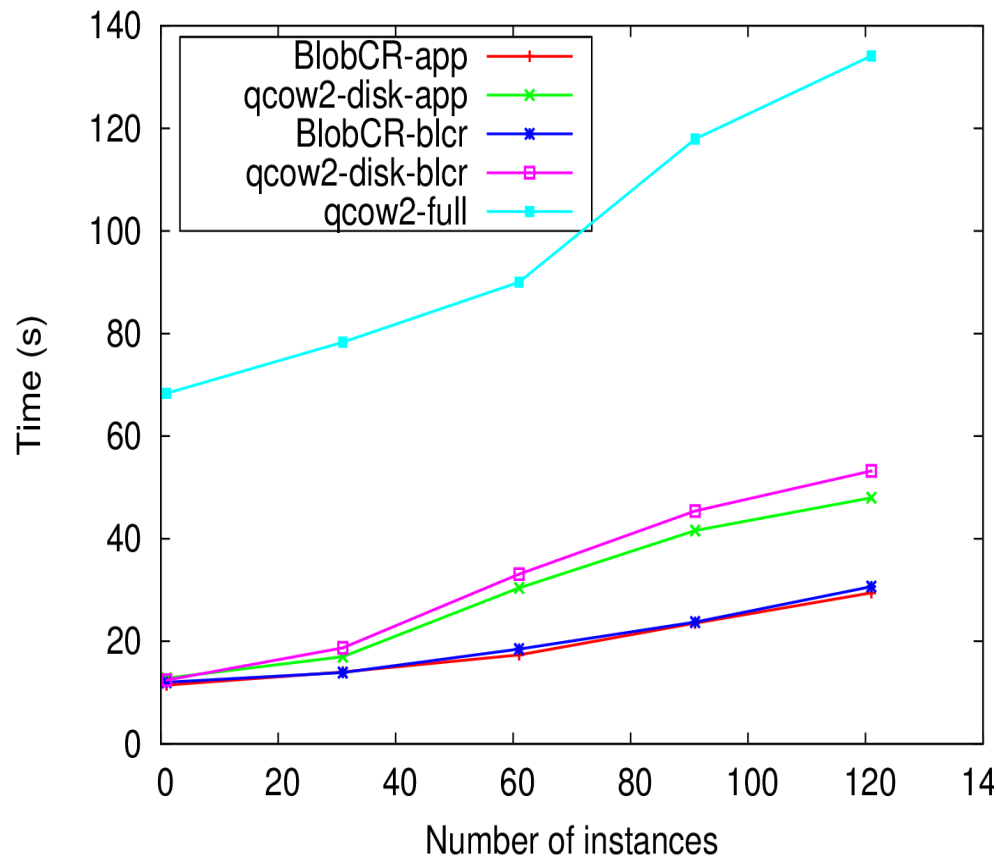
R/R, R/W, W/W



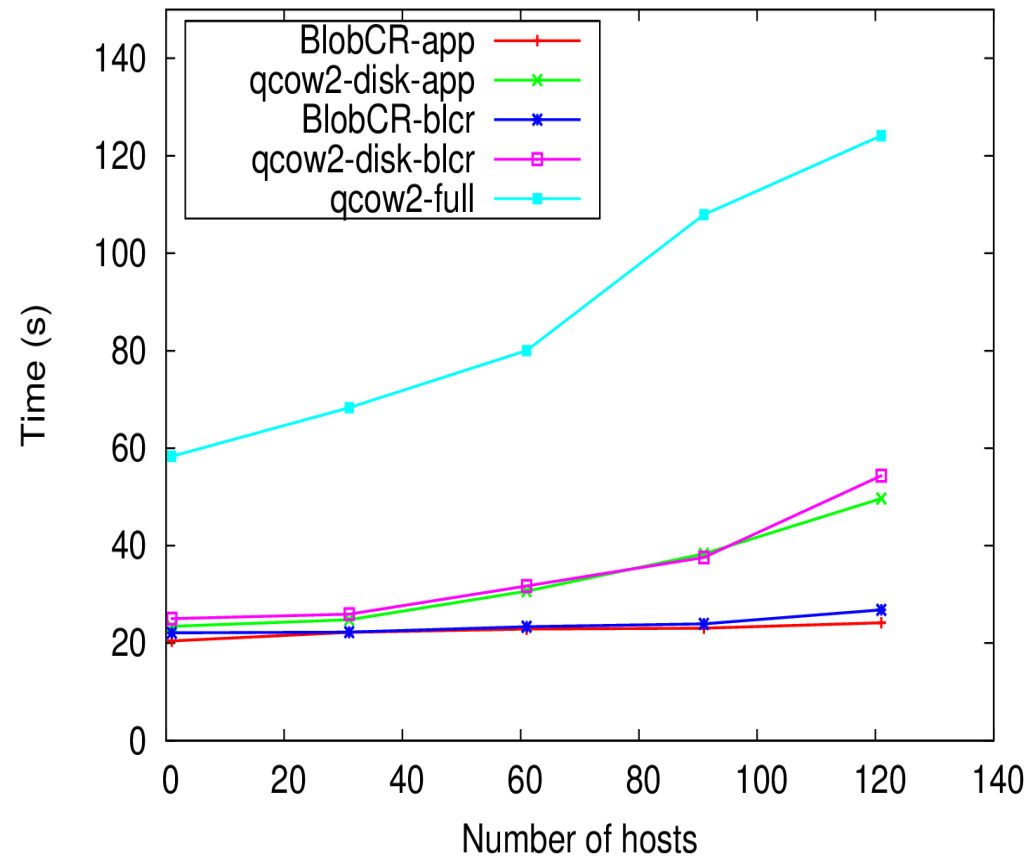
Using BlobSeer for cloud data management

- VM management to build a scalable, highly-available IaaS
 - BlobSeer internally used in the cloud for VM deployment and checkpointing
 - Integration in Nimbus
- Sharing application-level data in IaaS PaaS
 - Multiple VMs share application data through BlobSeer
 - BlobSeer exposes multiversioning to clients
 - Integrated within Nimbus, Azure
- Cost-effective storage service built on top of multiple clouds (sky computing)
 - BlobSeer relies on external, virtualized storage resources

Checkpoint/Restart performance



Speedup vs. PVFS + QCOW2:
8x VM level,
2x process + app level



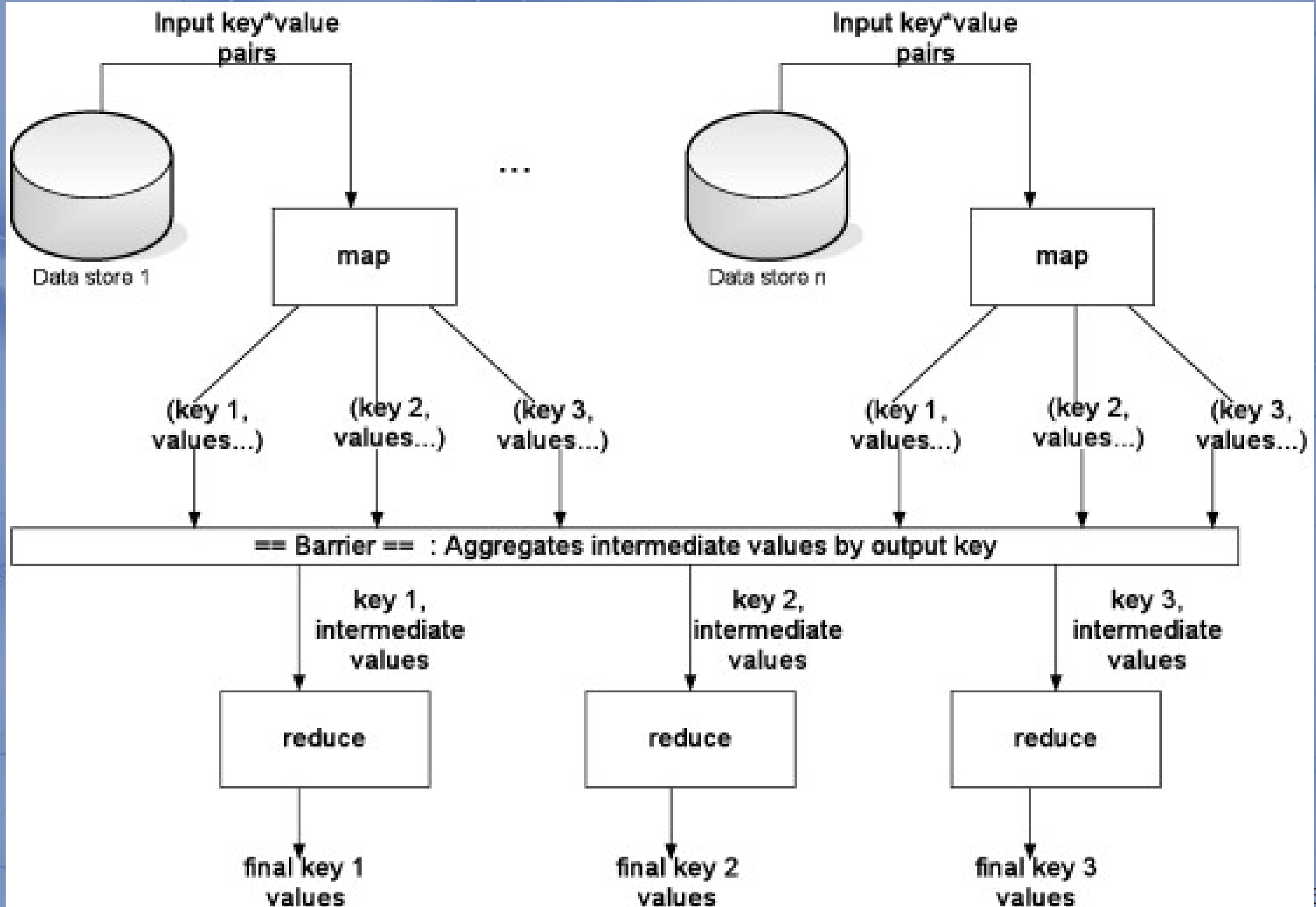
Speedup vs. PVFS + QCOW2:
6x VM level,
3x process + app level

#2: New programming models for very large-scale programming

- A **simple programming model** that applies to many **data-intensive** computing problems
- Approach: **hide messy details** within a runtime library
 - Automatic parallelization
 - Load balancing
 - Network and disk transfer optimization
 - Handling of machine failures
 - Robustness
 - Improvements to core library benefit all users of library!

Scientific contribution: MapReduce!

- Typical problem solved by MapReduce
 - Read a lot of data
 - **Map**: extract something you care about from each record
 - Shuffle and Sort
 - **Reduce**: aggregate, summarize, filter, or transform
 - Write the results
- Outline stays the same, Map and Reduce change to fit the problem
 - **map**(k, v) $\rightarrow \langle k', v' \rangle^*$
 - **reduce**($k', \langle v' \rangle^*$) $\rightarrow \langle k', v'' \rangle^*$



MapReduce: counting words

```
map(String input_key, String input_value):  
  // input_key: document name  
  // input_value: document contents  
  for each word w in input_value:  
    EmitIntermediate(w, "1");
```

```
reduce(String output_key,  
       Iterator intermediate_values):  
  // output_key: a word  
  // output_values: a list of counts  
  int result = 0;  
  for each v in intermediate_values:  
    result += ParseInt(v);  
  Emit(AsString(result));
```

Distributed grep
Distributed sort
Term-vector per host
Document clustering
Machine learning

Architecture and Scheduling

- One master, many workers
- Master assigns each map task to a free worker
- Master assigns each reduce task to a free worker
- Master detects worker failures
- Master notices particular input key/values that cause crashes in map(), and skips those values on re-execution

Sorting 1PB with MapReduce

November 22, 2008 at 1:55 AM



At Google we are fanatical about organizing the world's information. As a result, we spend a lot of time finding better ways to sort information using [MapReduce](#), a key component of our software infrastructure that allows us to run multiple processes simultaneously. MapReduce is a perfect solution for many of the computations we run daily, due in large part to its simplicity, applicability to a wide range of real-world computing tasks, and natural translation to highly scalable distributed implementations that harness the power of thousands of computers.

In our sorting [sort benchmark](#) sometimes you need to sort more than a terabyte, so we were curious to find out what happens when you sort more and gave one petabyte (PB) a try. One petabyte is a thousand terabytes, or, to put this amount in perspective, it is 12 times the amount of [archived web data](#) in the U.S. Library of Congress as of May 2008. In comparison, consider that the aggregate size of data processed by all instances of MapReduce at Google was on average 20PB per day in [January 2008](#).
benefits of va as an Olympi programs, we lessons useft help everyon

We are excited [System](#) as 10 computers in seconds on 9
It took six hours and two minutes to sort 1PB (10 trillion 100-byte records) on 4,000 computers. We're not aware of any other sorting experiment at this scale and are obviously very excited to be able to process so much data so quickly.

An interesting question came up while running experiments at such a scale: Where do you put 1PB of sorted data? We were writing it to 48,000 hard drives (we did not use the full capacity of these disks, though), and every time we ran our sort, at least one of our disks managed to break (this is not surprising at all given the duration of the test, the number of disks involved, and the expected lifetime of hard disks). To make sure we kept our sorted petabyte safe, we asked the Google File System to write three copies of each file to three different disks.

Hadoop: MapReduce for the masses



Cluster of machines running Hadoop at Yahoo! (Source: Yahoo!)

Word count example in Hadoop

```
public void map(WritableComparable key, Writable value, OutputCollector output, Reporter reporter) throws  
    IOException {
```

```
    String line = ((UTF8)value).toString();
```

```
    StringTokenizer itr = new StringTokenizer(line);
```

```
    while (itr.hasMoreTokens()) {  
        word.set(itr.nextToken());  
        output.collect(word, one);  
    }  
}
```

```
public void reduce(WritableComparable key, Iterator values, OutputCollector output, Reporter reporter) throws  
    IOException {
```

```
    int sum = 0;
```

```
    while (values.hasNext()) {  
        sum += ((IntWritable) values.next()).get();  
    }
```

```
    output.collect(key, new IntWritable(sum));  
}
```

Who uses Hadoop?

- Amazon/A9
- Facebook
- IBM: Blue Cloud?
- Joost
- Last.fm
- New York Times
- PowerSet
- Veoh
- Yahoo!

#3: New complexity models for very large-scale computing

Uncertainty			Scheduling		
Object	Nature	Origin	Type	Criterion	Problem ¹
computation duration	methodological	hardware, software	optimization	robustness	$R prec C_{\max}$
computation success	aleatory	hardware	evaluation	reliability	$R prec C_{\max}$
result correctness	epistemic	software, human	characterization	precision	$R online - time - nclv \sum C_i$

Uncertainty

Methodological: limitation(s) of the method (e.g., model simplification)

Epistemic: inaccessible knowledge (e.g., online task submission)

Aleatory: stochastic variability (e.g., hardware fault)

En guise de conclusion...



Today's challenge: Think Big!

