# iRODS:
# A Highly Customisable Data Management System To Face Big Data Challenges
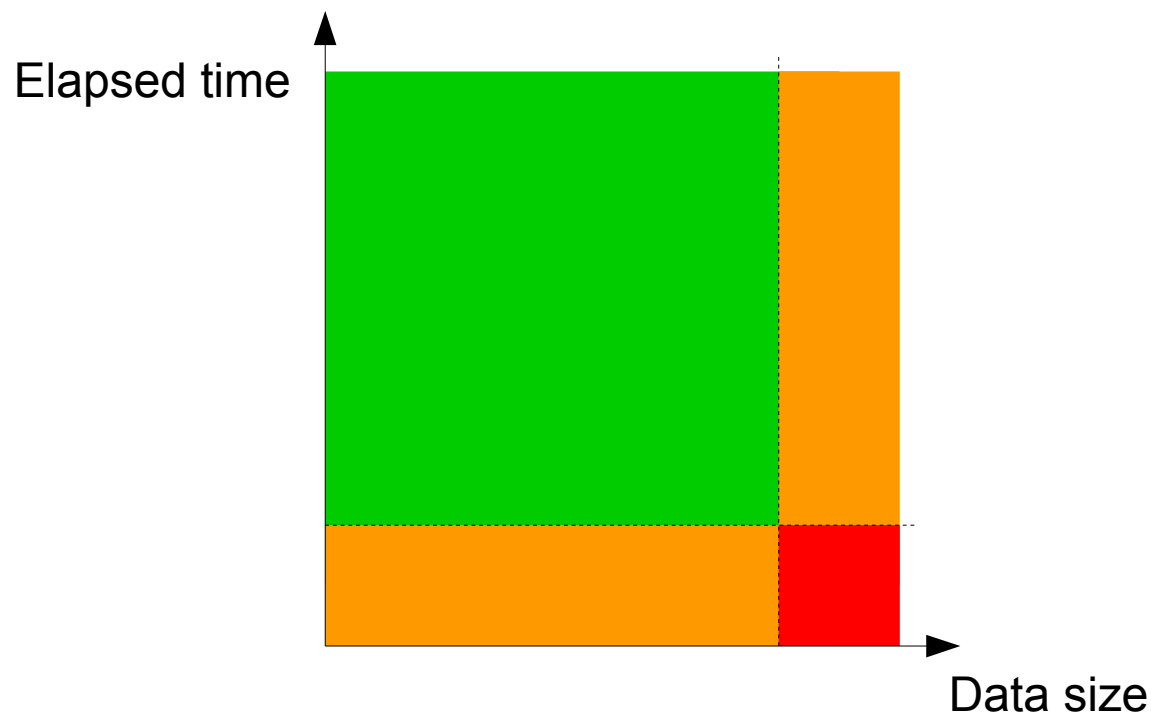
Formation iRODS – 13/02/2014

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

- Introduction to Big Data

- Data Management for Big Data

- Quick Overview of iRODS

- Use Cases

- France-Grilles Service Offering

## Big Data

*Big data usually includes data sets with sizes beyond the ability of commonly used software tools to capture, curate, manage, and process the data within a tolerable elapsed time.*[1]



Elapsed time

Data size

[1]C. Snijders, U. Matzat & U.-D. Reips: 'Big Data': Big gaps of knowledge in the field of Internet. International Journal of Internet Science 2012, 7, 1-5.

## Twitter's Challenges (Social Networks)

One of Twitter's challenges is to keep statistics of Tweets and Tweeted URLs:
- Several of them are retweeted by millions of followers
- At any time, a famous person can tweets a URL to millions of followers
- 143,199 Tweets per second record (3rd of August 2013)
- Top retweeted URLs is an important feature for many users
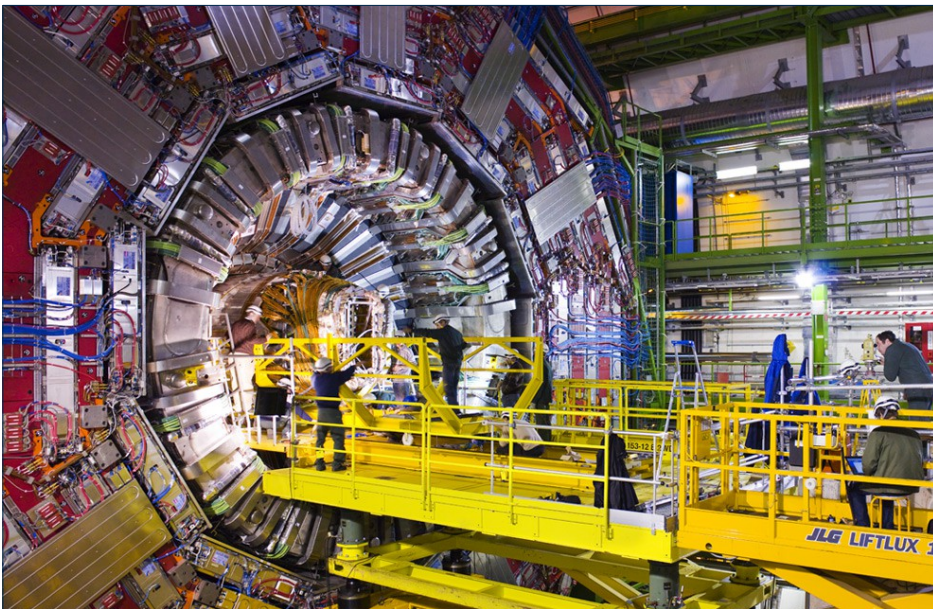


Key technologies:
- 250 millions tweets per day stored with MySQL
- Storm and Hadoop are used to proceed with unstructured and large dataset
- Cassandra is used for high velocity writes
- Vertica is used for analytics

https://blog.twitter.com/engineering

## LHC Data Analysis

The CERN is operating the Large Hadron Collider (LHC) in Geneva, Switzerland. This accelerator produced a huge amount of data:

- A 200-megapixel camera
- 40 million frames every second => 1000 TB/s
- Equipped with 4 detectors
- 27-kilometer circular collider
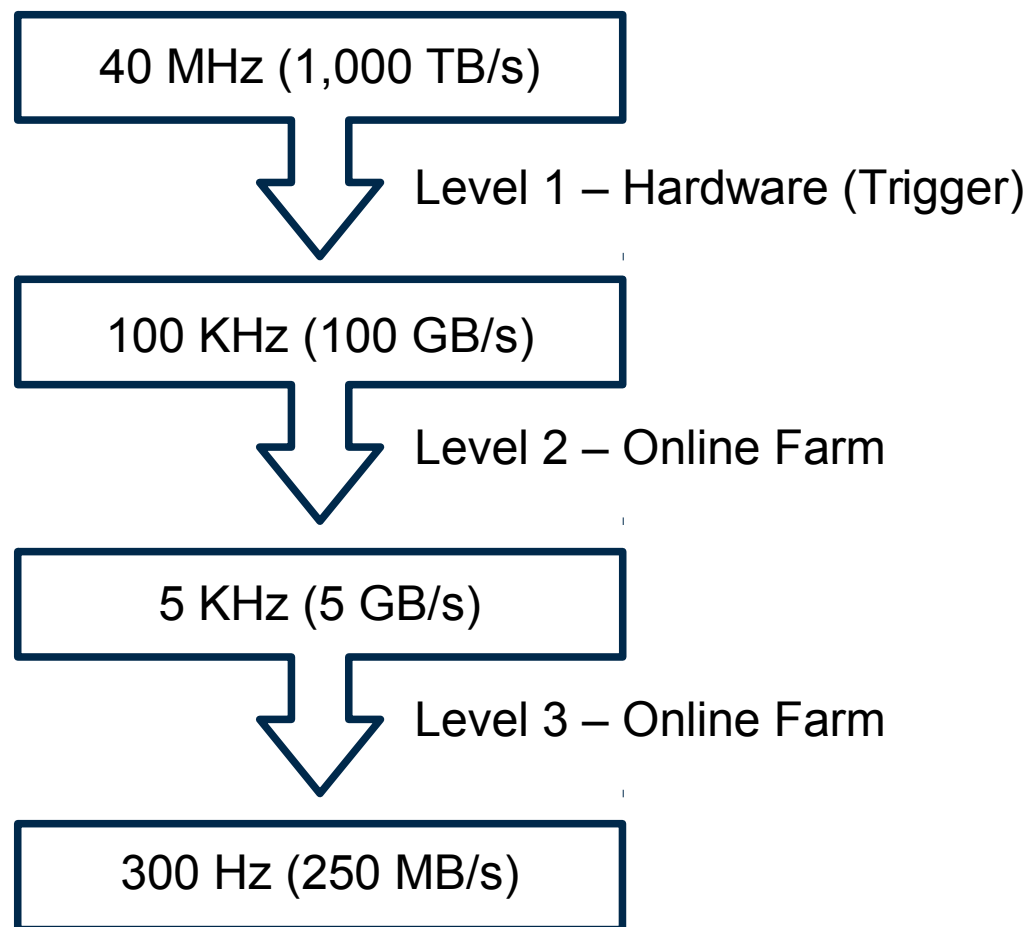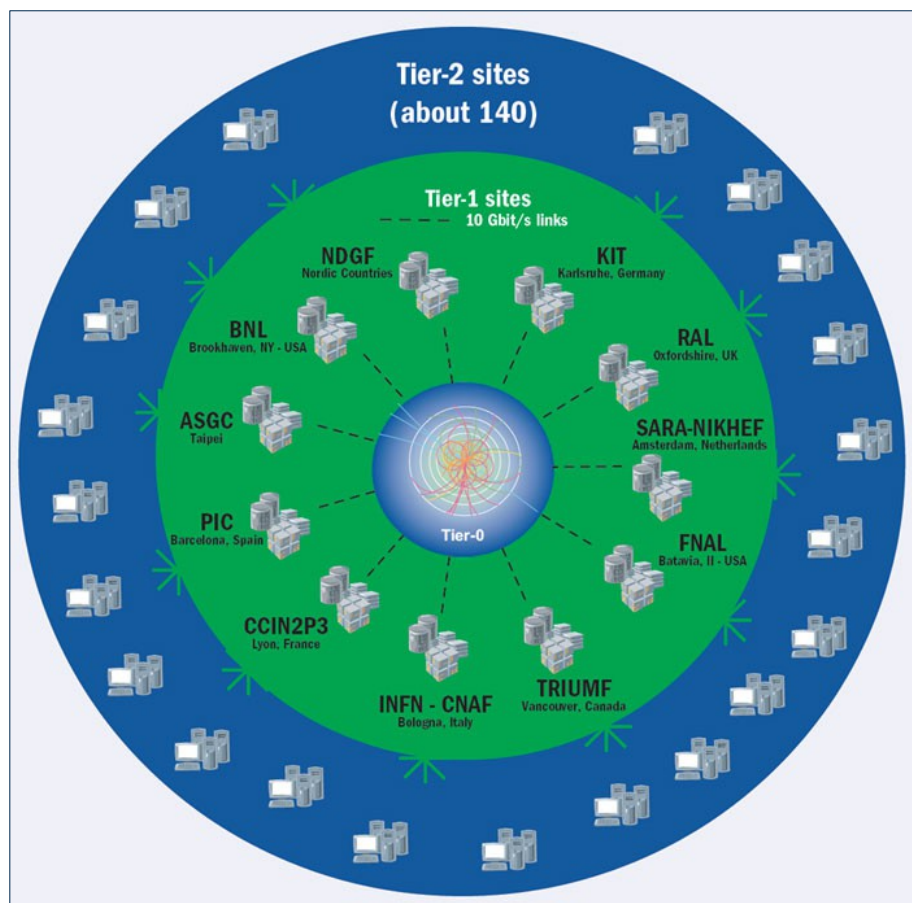- 25 PB of data per year (~3.14 x the height of Everest if all data were stored on CDs)



Key technologies:
- Grid Computing
- DPM and dCache storage technologies
- SRM, RFIO and XRootD protocols

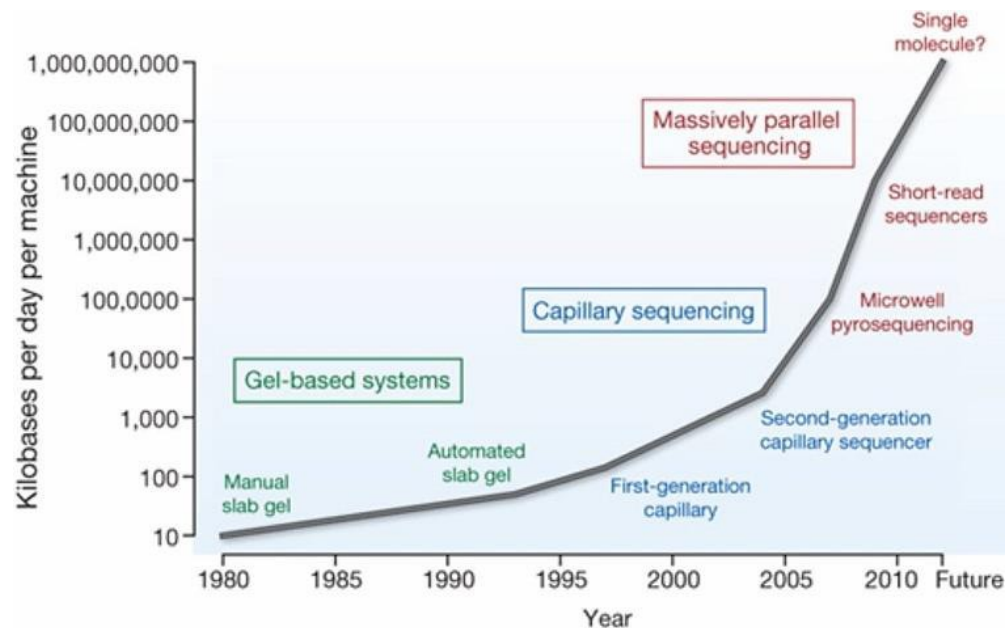# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

## LHC Data Analysis



40 MHz (1,000 TB/s)

Level 1 – Hardware (Trigger)

100 KHz (100 GB/s)

Level 2 – Online Farm

5 KHz (5 GB/s)

Level 3 – Online Farm

300 Hz (250 MB/s)

## DNA Sequencing



Stratton (2009) Nature

Illumina HiSeq2500:
- Up to 120GB per day → 44 TB per year
- Quick access for genome alignment
- Backup
- Typically sold with servers that can store data made in a year

## High-Throughput Imaging

Microscope:
- FEI Tecnai F30
- 16 million pixel camera
- 16-bit color depth
- Up to 40 frames per second
- ~ 1 GB/s $\rightarrow$ O(10) TB per day

IT Requirements:
- High-speed network (10Gb/s)
- High capacity and high performance storage system (hybrid solution)
- Backup (disk, tape)

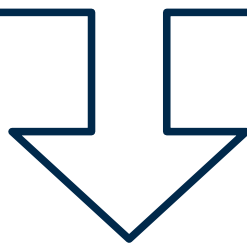## Data Management Context

Different scientific fields environment:
- Humanities and Social Sciences
- High Energy Physics
- Biology
- Biomedical Applications
- Astrophysics
- ...

**Various constraints, various needs for data management.**

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

- Introduction to Big Data
- Data Management for Big Data
- Quick Overview of iRODS
- Use Cases
- Available Infrastructures

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

## Big Data Software Landscape



http://www.bigdatalandscape.com/

## Infrastructure

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

- Data stored on sites in different locations
- Heterogeneous storage:
  - Data format: flat files, databases, data stream,...
  - Storage media, server hardware
  - Data access protocols, information systems
- Heterogeneous OS on both clients and servers side

One Soft to rule them all, One Soft to find them,

One Soft to bring them all and in the data center bind them

**i·R·O·D·S**

Integrated Rule-Oriented Data System

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

- Introduction to Big Data
- Data Management for Big Data
- Quick Overview of iRODS
- Use Cases
- Available Infrastructures

## What is iRODS ?

http://www.irods.org/index.php

iRODS (iRule Oriented Data Systems) is a data grid software system providing a transparent access to data spread over different physical locations and heterogeneous storage technologies:

- Project has been started in 2006 by the DICE team (UNC, San Diego)
- Open Source
- Financed by NSF and NARA



User

Search, Access, Add and
Manage Data / Metadata

iRODS Data Server
(disk, tape, etc)

iRODS Rule Engine
(manage policies)

iRODS Metadata Catalog
(manage metadata)

**iRODS Data System**

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

## What is iRODS ?

- Virtualised storage servers in a *Zone* (administrative domain):
  - One or several servers connected to a centralised Metacatalogue → logical view of the data in a given zone
  - Data servers can be spread geographically within one zone
  - Possibility to have different zones interconnected

- Data management policies expressed with rules:
  - Can be triggered automatically for various actions (put, get, …)
  - Can be run manually
  - Can be run in batch mode

- Client interactions with iRODS:
  - APIs (C, Java, PHP, Python), shell commands, GUIs, web interfaces

- Further informations:
  - → `http://storageconference.org/2013/Presentations/Moore.pdf`

## E-iRODS: Enterprise Quality Data Management



- Objectives:
  - To ensure the long-term sustainability of iRODS
  - To provide a fully tested software by using complementary process of testing, packaging, and expertise developed at RENCI
  - To provide binary packages
  - To release E-iRODS as an Open Source software

- For more informations:
  - → `http://eirods.org`

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

- Introduction to Big Data

- Data Management for Big Data

- Quick Overview of iRODS

- Use Cases

- Available Infrastructures

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

## iRODS at CC-IN2P3



- iRODS is a key service
- In production since 2008
- User support provided by several experts
- Used by several projects (Adonis, BaBar, biology, biomedical apps, …)
- 5 PB in 2012 managed with iRODS
- Connected to tape library through HPSS driver
- Replication with Paris, Grenoble and Montpellier (CINES)
- Customised with several rules

- For more informations:
    - → `http://cc.in2p3.fr/IRODS,2059`
    - → `http://storageconference.org/2012/Presentations/M14.Nief.pdf`

## Rule examples: biomedical data



- Human and animal data (fMRI, PET, MEG, …)
- Usually in DICOM format
- Need to anonymised human data
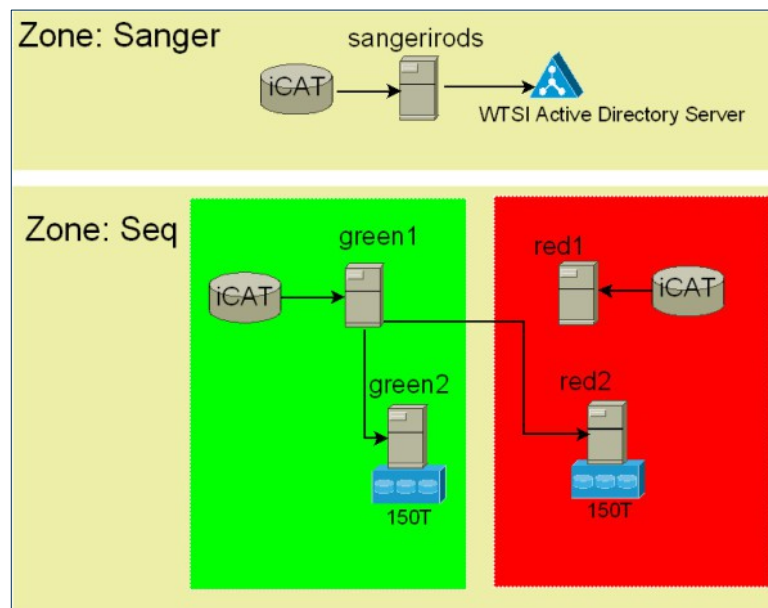- Need to do metadata search on DICOM files

**Rule Engine**

- Check for anonymisation of the file: send a warning if not true
- Extract a subset of metadata (based on a list stored in iRODS) from DICOM files
- Add these metadata as user defined metadata in iRODS

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

## Genomic Data Management with iRODS

WTSI Use Case:[1]
- Managing and accessing sequencing Binary Alignment/Map (BAM) files
- 500 TB SAN Storage
- Integrated in the sequencing pipeline
- Fine-grained access control
- Data replication
- Metadata on alignement are automatically added
- Data federation with other research institutes



[1]G.-T. Chiang, P. Clapham, G. Qi, K. Sale & G. Coates: Implementing a genomics data management system using iRODS in the Wellcome Trust Sanger Institute. BMC Bioinformatics 2011, 12, 361.

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

## SILS Lifetime Library

Objectives:
- Provide trustworthy and easy to use services that help students and alumni to sustain, extend, and use the information resources that compose their knowledge base over a lifetime
- Solution independent from device (laptop, desktop, mobile phones, …)
- Serve as a link to alums who stay in touch and participate in campus activities
- Integrate a 120 PB infrastructure based on cloud services

Achievements:
- Development of the iDrop (iRODS GUI)
- Infrastructure is made available through a web portal
-  Distributed mass storage arrays is integrated using the iRODS middleware

- For more informations:
  - → `http://lifetime-library.ils.unc.edu`

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

## Other examples

- Astrophysics: Auger supernova search
- Atmospheric science: NASA Langley Atmospheric Sciences Center
- Biology: Phylogenetics at CC IN2P3
- Climate: NOAA National Climatic Data Center
- Cognitive Science: Temporal Dynamics of Learning Center
- Computer Science: GENI experimental network
- Cosmic Ray: AMS experiment on the International Space Station
- Dark Matter Physics: Edelweiss II
- Digital Library French National Library, Texas Digital Libraries
- Earth Science: NASA Center for Climate Simulations, Vhub - vulcanism
- Ecology: CEED Caveat Emptor Ecological Data
- Engineering: CIBER-U
- High Energy Physics: BaBar
- Hydrology: Institute for the Environment, UNC-CH; Hydroshare
- Genomics: Broad Institute, Wellcome Trust Sanger Institute, NGS
- Indexing: Cheshire
- Institutional repository: Carolina Digital Repository
- Medicine: Sick Kids Hospital
- Neuroscience: International Neuroinformatics Coordinating Facility
- Neutrino Physics: T2K and dChooz neutrino experiments
- Oceanography: Ocean Observatories Initiative
- Optical Astronomy: National Optical Astronomy Observatory

# iRODS: A Highly Customisable Data Management System To Face Big Data Challenges

- Introduction to Big Data

- Data Management for Big Data

- Quick Overview of iRODS

- Use Cases

- France-Grilles Service Offering