# Clouds in biomedical sciences
# Part IV – entering a new world

Vincent Breton

July 28th 2014

Enrico Fermi school of physics

# Session IV: the future

- Welcome to a new world

- Learn from history to prepare future: an introduction to Big Data

- What I do of my spare time…

# A new world beyond the standard model

- For more than 30 years, validation of the standard model
    - Electroweak physics at LEP
    - Top quark discovery at TEVATRON
    - Higgs Boson discovery at LHC
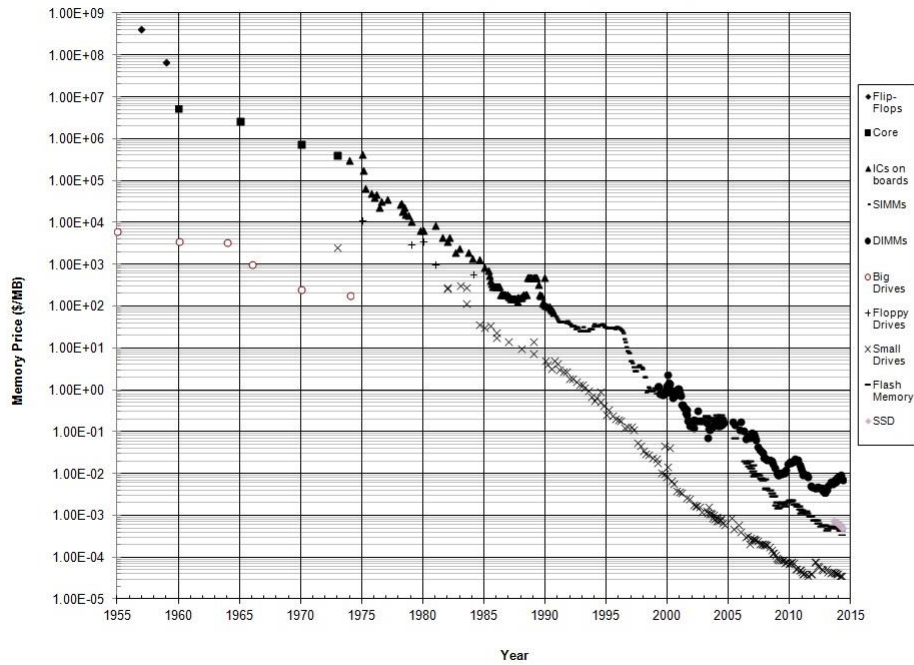- New exploratory phase beyond the standard model
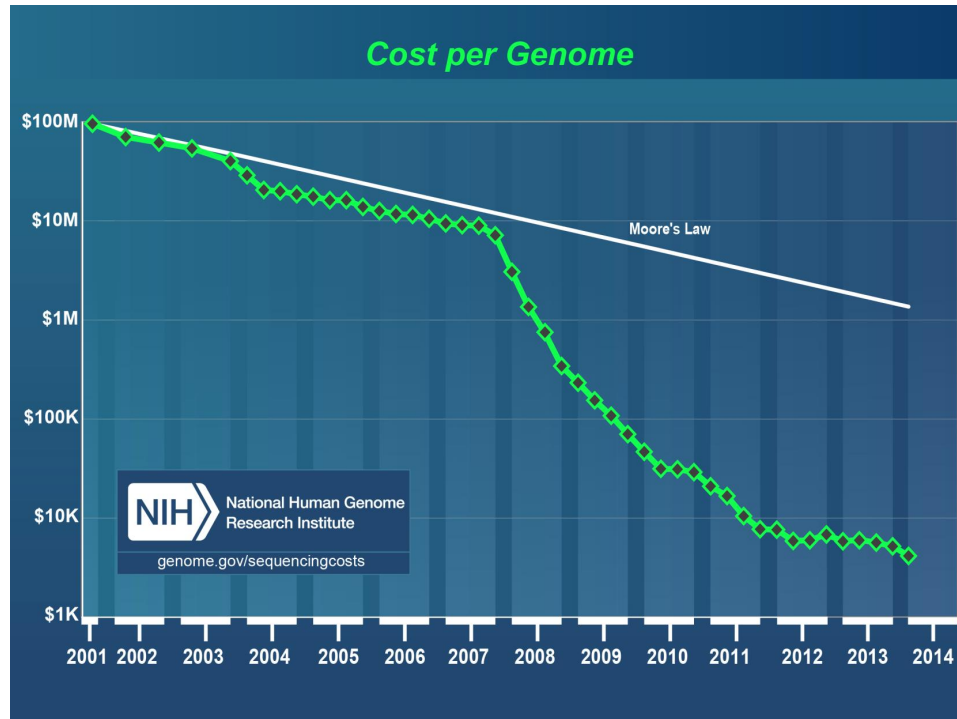    - Where is the new physics?

# A new world without Moore's law

- Moore's law does not apply any more to storage capacities… nor to sequencing data production



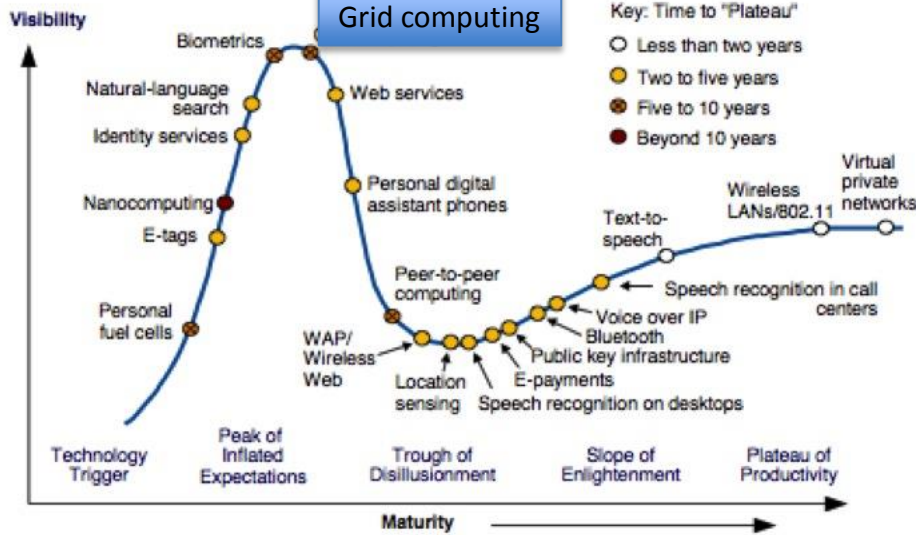Historical Cost of Computer Memory and Storage



Cost per Genome
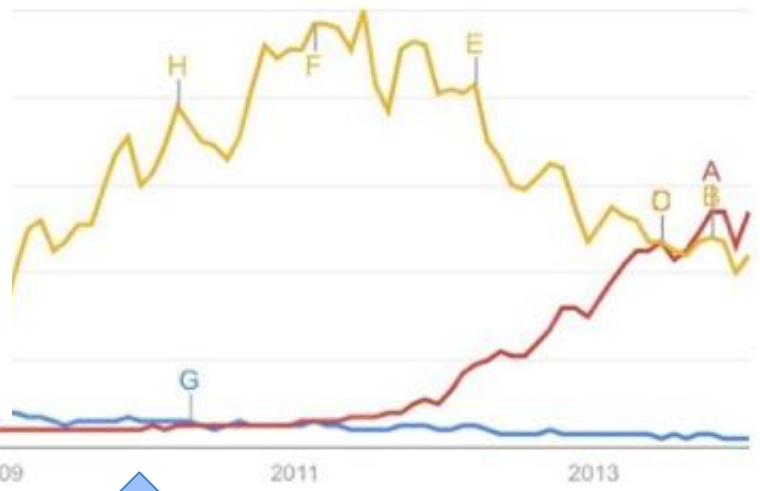
# It takes many years from hype to production quality



Gartner Emerging Technologies Hype Cycle 2002

Grid computing

A long way to cloud maturity
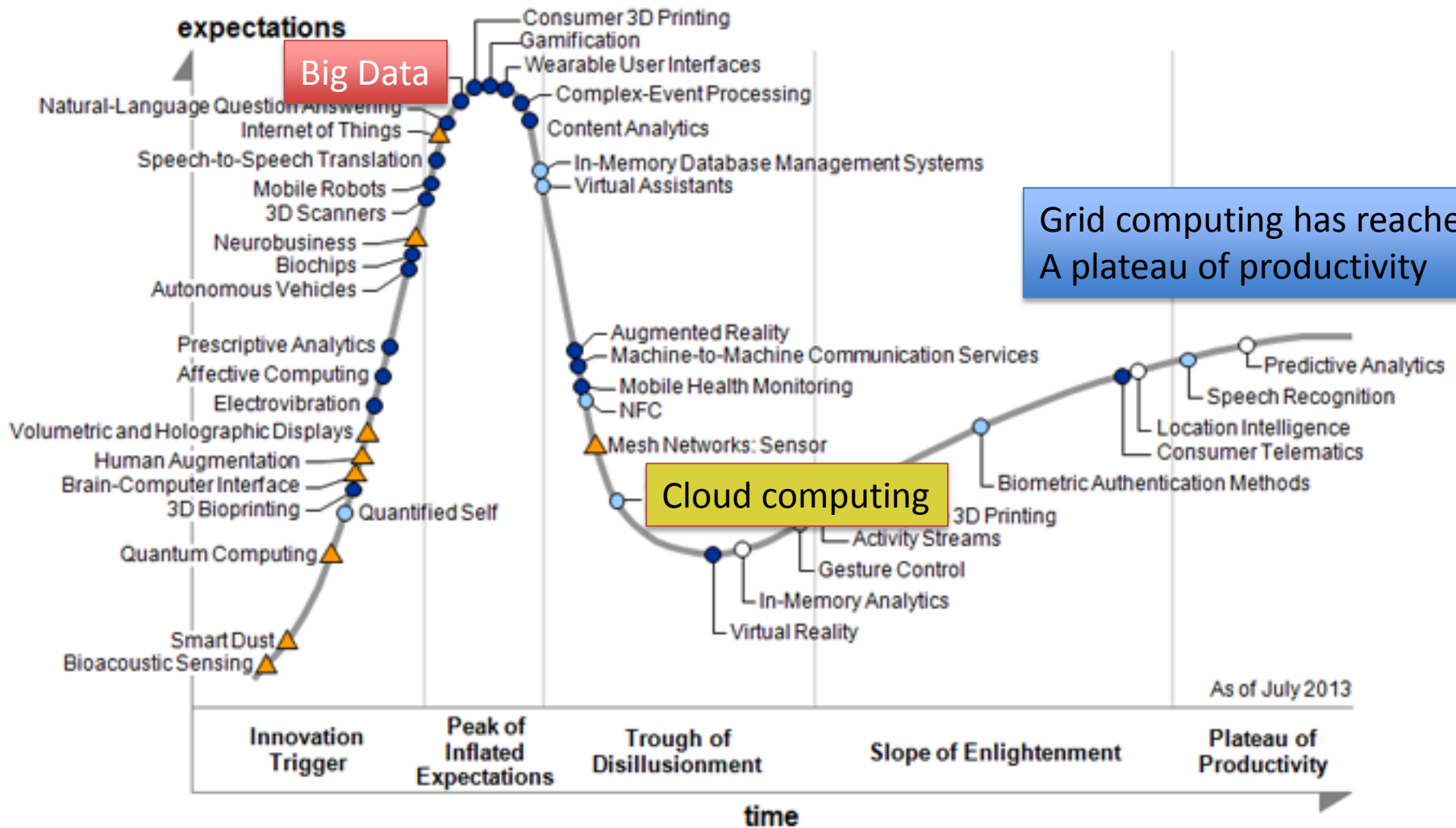
Grid peak of expectations back in 2002

Grid maturity

# Gardner hype curve for 2013
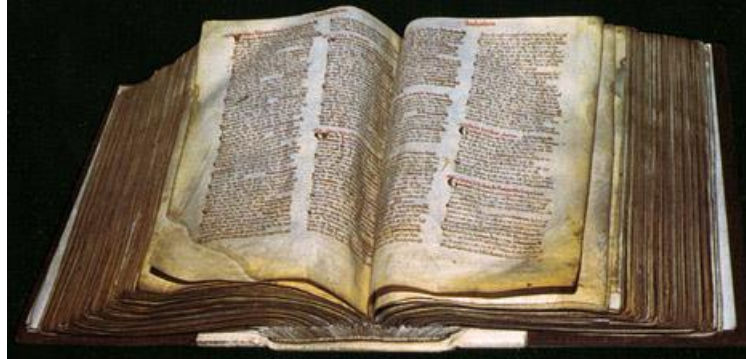
# Learning from history to build the future

- The greatest achievement of grids is not the capacity it has built
  - Obsolescence in three years for the hardware
  - Obsolescence of the grid middleware
- The greatest achievement are the human networks it has created
  - Fantastic human adventure

# Learning from history: the Domesday Book (1087)



- Manuscript record of the great survey, completed in 1086 on orders of [William the Conqueror](#)

*«While spending the Christmas time of 1085 in Gloucester, William had deep speech with his counsellors and sent men all over England to each shire to find out what or how much each landholder had in land and livestock, and what it was worth»*                                    *Anglo-Saxon chronicle*

- Absolute authority to define property rights since Middle Age

*for as the sentence of that strict and terrible last account cannot be evaded by any skilful subterfuge, so when this book is appealed to ... its sentence cannot be quashed or set aside with impunity. That is why we have called the book 'the Book of Judgement' ... because its decisions, like those of the Last Judgement, are unalterable.*                                    *Richard Fitzneal, Dialogus de Scaccario, 1179*

# Big data issues (I/II)

- **Data collection**
  - Every shire visited by a group of royal officers (1085-1086)
  - The unit of inquiry was the Hundred (a subdivision of the county)

- **Data veracity**
  - return for each Hundred was sworn to by twelve local jurors, half of them English and half of them Normans.

- **Data analysis**
  - names of the new holders of lands and assessments on which their tax was to be paid
  - national valuation list, estimating the annual value of all the land in the country

# Big Data issues (II/II)

- **Data presentation**
  - Properties listed by fiefs
  - Properties listed by owner categories
    - king's holdings
    - holdings of churchmen and religious houses
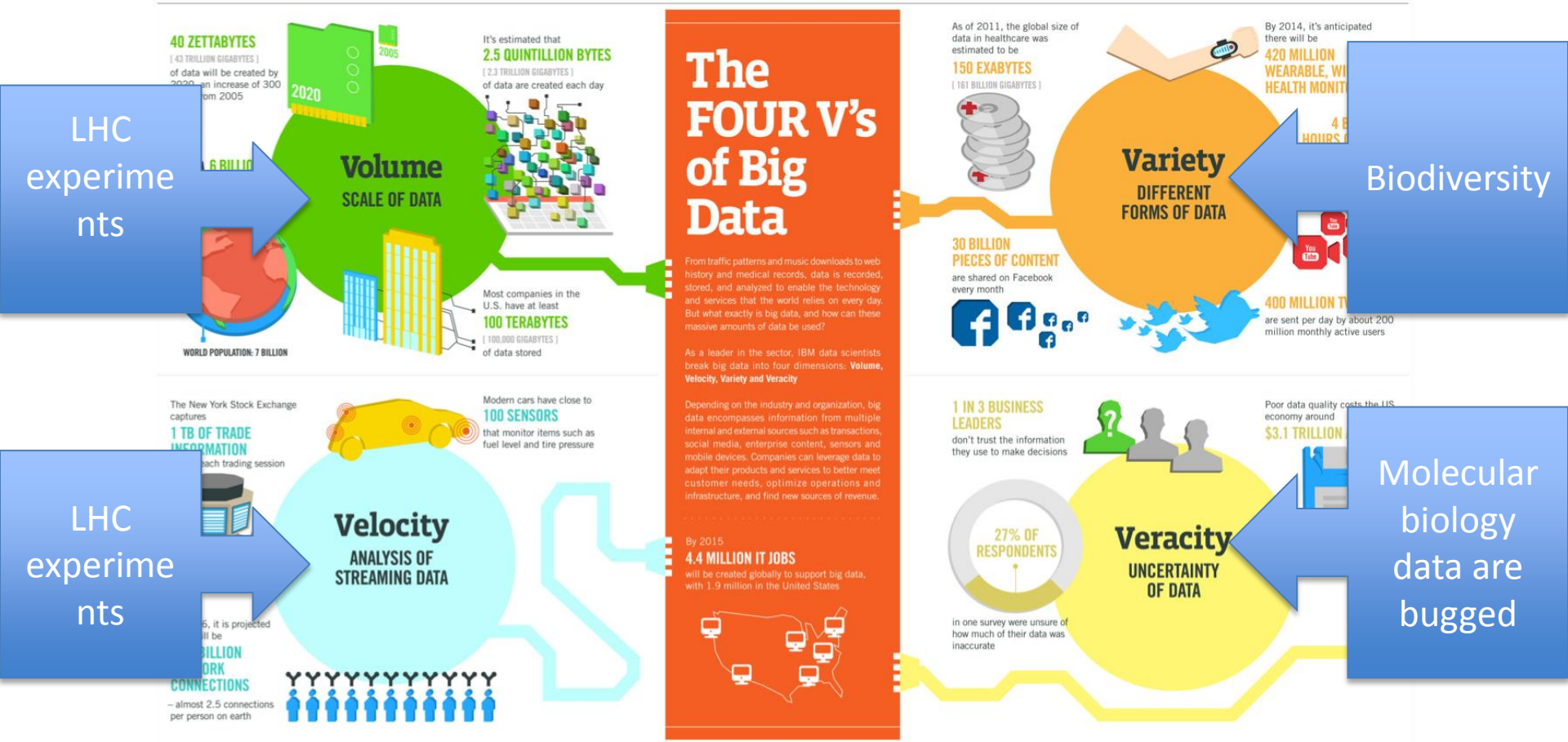    - Aristocrats
    - Lay men

- **Data preservation**
  - Preservation in the Royal Treasury in Westminster till 19th century
  - Stored at UK National Archives in Kew
  - 1986: digital version
  - 2002: access problem to digital version



TERRA ARCHIEPI CANTVAR. *In Waletone hvnd.*
II. Archieps Lanfranc ten in dnio *Croindene*. T.R.E. se
defd p quat xx . hid . 7 modo p . xvi . hid 7 una v. Tra . e̅
. xx . car . In dnio sunt . iiii . car . 7 xlviii . uilli 7 xxv . bord . cu̅
. xxxiiii . car . Ibi æccla . 7 un molin de . v . sol . 7 viii . ac pti . Silua:
de . cc . porc.
De tra huj ᴔ ten Restold vii . hid de archiepo. Radulf . i . hida.
7 inde hnt . vii . lib 7 viii . sol . de gablo.
Totu̅ T.R.E 7 post: ual . xii . lib. Modo: xxvii. lib archiepo.
Hominib; ej . x . lib 7 x . solid.
Ipse archieps ten *Ceiha* de uictu monachoz . T.R.E. se defd
p . xx . hid . 7 m p . iiii . hid . Tra . e̅ . xiiii . car . In dnio sunt . ii . car.
7 xxv . uilli 7 xii . cot . cu̅ . xv . car . Ibi æccla 7 v . serui . 7 una
ac pti . Silua: de . xxv . porc.
T.R.E. 7 post: ualuit . viii . lib. Modo: xiiii . lib. *In Brixiestan hd.*

# Big Data 4 Vs



The infographic labeled "The FOUR V's of Big Data" with four sections: Volume (Scale of Data), Velocity (Analysis of Streaming Data), Variety (Different Forms of Data), and Veracity (Uncertainty of Data). Blue arrow callouts point to each: "LHC experiments" → Volume, "LHC experiments" → Velocity, "Biodiversity" → Variety, "Molecular biology data are bugged" → Veracity.
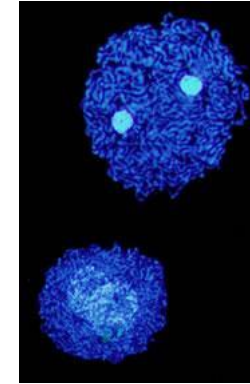
# Data volumes: the example of metagenomics

*Metagenomics* is the study of genetic material recovered directly from environmental samples.



Smallest non viral genome: *Carsonella ruddii* (0,16Mbp)

### Evolution of sequencing techniques

| | |
|---|---|
| Sanger technology | 500 base pairs (bp) |
| 454 technology | $10^5$ 400-600 bp reads |
| Illumina Technology | $10^6$ 100 bp reads |
| TARA project | $10^7$ 100-400 bp reads |





Largest genome: *Polychaos dubium* (670Gbp)

Tara @ http://oceans.taraexpeditions.org/

# Cost per Genome is decreasing faster than Moore's law



**Cost per Genome**

Moore's Law

National Human Genome Research Institute

genome.gov/sequencingcosts

# Consequence: over 2500 Next Generation Sequencing machines in 900+ research centers in the world



> 60PB per year

Source: omicspmaps.com

# Welcome to Auvergne, at the heart of France

1,35 Million inhabitants
26013 km$^2$

# Auvergne at the heart of Uranium production in France

Map of uranium mines in metropolitan France

1949: first attempt to extract uranium ore in France in Lachaux (Auvergne)

In 50 years:

- 53 Million tons extracted in France till 2001
- 76000 tons of uranium ore produced in > 200 mines

Départements concernés
- Sondages et tranchées
- Travaux miniers souterrains (TMS) ou Mines à Ciel Ouvert (MCO)

**Production d'uranium (en tonnes)**
- ≃ 0 t
- > 0 t - 50 t
- > 50 t - 100 t
- > 100 t - 500 t
- > 500 t - 1 000 t
- > 1 000 t - 10 000 t
- > 10 000 t
- Information manquante

Le code couleur renseigne sur la masse d'uranium métal produite à partir du minerai extrait des mines concernées (et non pas sur le tonnage du minerai).

En France, pour produire 1 tonne d'uranium, il a fallu extraire, en moyenne, 1400 tonnes de minerai (stériles uranifères non compris).

# ZATU, a Long Term Ecological Research dedicated to life under natural ionizing radiation


Natural radioactivity


Storage sites of uranium ore extraction residues

- Society in uranium rich territories
  - Social impact of uranium extraction
  - Preserving the long term memory
- Characterization, behavior and transfer of radionucleids
  - long term future of radionucleids in storage sites
- Impact of radiation on living systems
  - Multigenerational effects of chronic exposure to radiation

- From the Chernobyl environment, a coherent picture of predictable radiation-induced effects for low-dose-rate exposures has not emerged
  - Contradictory experimental evidences from Chernobyl exclusion zone
- Need to collect more data from Chernobyl exclusion zone but also from other ecosystems under chronic low dose exposure
  - Radioactive water sources
- Point 0: what happens in "total" absence of radioactivity?



Photographs of abnormalities in barn swallows. (a) Normal phenotype. (b–d) Partially albinistic plumage. (e) and (f) Deformed beak. (g) Deformed air sacks. (h) and (i) Bent tail feathers.

# ZATU strategy

Multidisciplinary long term observation of selected sites in Auvergne, Massif Central and Massif Armoricain

- Radionucleid chemical speciation
- Industrial heritage
- Biodiversity survey

**Characterization**

**Transfer**

- Radionucleid migration
- Interaction of radiation with living organisms
- Territory administration and responsabilities

- Interactions and retroactions between matter and living systems
- Risk evaluation
- Prevention tools

**Environmental impact**

Significant production of scientific data (geography, ecology, biology, metagenomics, chemistry, physics, social sciences)

How to make all these data speak to each other is a huge challenge

# Conclusion

- Grid computing has allowed building a truly multidisciplinary distributed IT infrastructure
  - Greatest achievement: human networks
- Cloud computing allows extending the grid functionalities
  - All sciences will benefit even more
  - Still a long way to the plateau of maturity
  - Scientific gateways and pilot agent platforms allow a smooth transition from grids to clouds
- Big Data is the next frontier
  - Volume will not be necessarily the most difficult challenge

# Which data produced today will still be used in 900 years?