## **ESFRI/e-IRG Collaboration**

#### **Yannis Ioannidis**

#### MaDgIK Lab University of Athens & ATHENA Research Center







Ερευνητικό Κέντρο Αθηνά

Ερευνητικό Κέντρο Καινοτομίας στις Τεχνολογίες της Πληροφορίας, των Επικοινωνιών, της Γνώσης

#### **ESFRI: European Strategy Forum on Research Infrastructures**

- Set up by the EU Council of Research Ministers in 2002
- Brings together representatives of Ministers of the 28 Member States, 9 Associated States, and of the EC
- Supports a coherent and strategy-led approach to policy making on Research Infrastructures
- Mandate to develop a Roadmap (2006) and its updates (2008, 2010)
- Evaluation and new Roadmap in the plans (... 2016)

## **ESFRI – The Roadmap**

New pan-European RIs or major upgrades to existing Ris



- Needs of European research communities in the next 10-20 years
- All fields of Sciences and Technologies, regardless of possible location
- Published: 2006, updated: 2008 & 2010
  - 48 projects
  - Financial investment: ~20 b€, long term commitment for operation: ~2 b€/year

## The ESFRI Process (1)



## The ESFRI Process (2)



## **ESFRI Roadmap 2010**

Social Sc. & Hum. (5)	Life Sciences (13)		Environmental Sciences (9)		Energy (7)	Material and Analytical Facilities (6)	Physics and Astronomy (10)		e-Infra- structure s (1)
SHARE	BBMRI	ELIXIR	ICOS	EURO- ARGO	ECCSEL	EUROFEL	ELI	TIARA	PRACE
European Social Survey	ECRIN	INFRA FRONTIER	LIFEWATCH	IAGOS	Windscanner	EMFL	SPIRAL2	СТА	
CESSDA	INSTRUCT	EATRIS	EMSO	EPOS	EU-SOLARIS	European XFEL	E-ELT	SKA	
CLARIN	EU- OPENSCREEN	EMBRC	SIAEOS	EISCAT_3D	JHR	ESRF Upgrade	KM3NeT	FAIR	
DARIAH	Euro BioImaging	ERINHA BSL4 Lab		COPAL	IFMIF	NEUTRON ESS	SLHC-PP	ILC- HIGRADE	
	ISBE	MIRRI			HiPER	ILL20/20 Upgrade			
	ANAEE				MYRRHA				
								Distributed research infrastructures	
								Single sited research infrastructures	



## Many <u>ARE</u> e-Infras

- CLARIN: Language resources ... enabling eHumanities
- DARIAH: Digital RI for the Arts & Humanities
- LIFEWATCH: An e-science and tech infrastructure for biodiversity and ecosystem research
- ELIXIR: Platform for biological data
- PRACE: HPC

. . .

## All need e-Infras

. . .

- EURO-ARGO: Ocean floor not only instruments ... but associated data streams and data centers
- IFMIF: Generation of a materials irradiation database ... for fusion reactors
- BBMRI: Bio-banks, bio-molecular resource centers ... and bio-computational tools

#### **ELIXIR** Safeguarding the results of life science research in Europe

#### **ELIXIR's mission**

To build a sustainable European infrastructure for biological information, supporting life science research and its translation to:

environment

medicine

bioindustries

society

# A distributed infrastructure to scale with the challenge

NIR NODE

ELIXIR data infrastructure for Europe's life science research sector

ELIXIR Nodes build local bioinformatics capacity throughout Europe

ELIXIR Nodes build on national strengths and priorities



## What is bioinformatics?

- The science of storing, retrieving and analysing large amounts of biological information
- An interdisciplinary science involving biologists, biochemists, computer scientists and mathematicians



• At the heart of modern biology

#### New building for ELIXIR Hub now open











#### Data resources in life science

- Diverse
- Many
- Disperse

# ~1800

#### molecular biology data resources

- Genomics Databases (non-vertebrate) (17.9%)
- Protein sequence databases (12.9%)
- Human Genes and Diseases (9.8%)
- Structure Databases (9.7%)
- Metabolic and Signaling Pathways (9.3%)
- Nucleotide Sequence Databases (8.8%)
- Human and other Vertebrate Genomes (7.1%)
- Plant databases (7.1%)
- RNA sequence databases (4.9%)
- Microarray and other Gene Expression Databases (4.5%)
- Other Molecular Biology Databases (3.3%)
- Immunological databases (1.8%)
- Organelle databases (1.6%)
- Proteomics Resources (1.2%)
- Cell biology (0.2%)

#### Nucleic Acids Research

Oxford Journals > Life Sciences > Nucleic Acids Research > Database Summary Paper Categories

#### 2012 NAR Database Summary Paper Category List

Nucleotide Sequence Databases RNA sequence databases Protein sequence databases Structure Databases Genomics Databases (non-vertebrate) Metabolic and Signaling Pathways Human and otabases (non-vertebrate) Human And other Vertebrate Genomes Human Genes and Diseases Microarray Data and other Gene Expression Databases Proteomics Resources Other Molecular Biology Databases Organelie databases Plant databases Plant databases Cell biology

Compilation Paper
 Category List
 Alphabetical List
 Category/Paper List
 Search Summary Papers

Online ISSN 1362-4962 - Print ISSN 0305-1048

Oxford University Press is not responsible for the content of external internet sites

Copyright © 2012 Oxford Journals

Compilation Paper
 Category List
 Alphabetical List
 Category/Paper List

Nucleic Acids Research annual Database Issue and the NAR online Molecular Biology Database Collection in 2012. MY Galperin, GR Cochrane – Nucleic Acids Research, 2011

#### **Industry and Innovation**

- Challenge to access and understand services
  - EU wide engagement programme - Build on strong local interactions through Nodes!
- Small companies Big Data
  - ELIXIR Cloud for SMEs



# A Life-science infrastructure of interoperable and integrated integrated data-services Will the improvement we see in an LPS-challenged mouse translate to a measurable clinical benefit?

Data interoperability, vocabulary and ontology services

user data	Europe PMC	EGA	Expression Atlas	
	НРА	Gene Ontology	UniProt	

# Future challenges for life-science data services

Scale and Sustain funding Managing and interoperate big and heterogeneous data

Capacity Compute. Capability Storage

Integrating clinical and translational data

Privacy and ethical concerns

Distributed infrastructure with >1M users

Algorithms to data – clouds, research environments...

### The data deluge

- Computer speed and storage capacity is doubling every 18 months and this rate is steady
- DNA sequence data is doubling every 6-8 months over the last 3 years and looks to continue for this decade

#### DATA EXPLOSION

The amount of genetic sequencing data stored at the European Bioinformatics Institute takes less than a year to double in size.



Source: *Nature News* & *Comment*, June 2013

#### Future scale of genomic data

- 1 Human Genome is 10<sup>9</sup> bases
- Human genome for all UK citizens 65m x 10<sup>9</sup> = 6x10<sup>16</sup>
- SNPs for all UK citizens  $65m \times 10^6 = 6\times 10^{13}$
- 1 petabyte = 10<sup>15</sup>
- EBI currently has ~16 petabytes of storage

#### Sequencing every child born in the EU?

5 million births per year

3 petabytes of raw data every day 9 petabases of DNA every week

Storing only variants: much more feasible

#### Clinical Data Raw, Anonymised and Summary Data



MAD IL

#### The scientific reason for ELIXIR

- Data is an essential commodity for life-science research
- Ten years ago, finding the connection between a gene and a characteristic (e.g., drought tolerance, risk of heart disease) could take years; now it takes minutes



Image courtesy of Genome Research Ltd.

- Data analysis is now the bottleneck in life-science research
- ELIXIR is our only realistic hope of easing that bottleneck

#### **Benefits of ELIXIR**

**ELIXIR** will contribute to European innovation by:

- Optimising access to and exploitation of life-science data
- Ensuring longevity of the data, thereby protecting investments already made in research
- Enhancing the quality of European research by supporting national efforts to increase the competence and number of bioinformatics users through training
- Strengthening the global position and influence of Europe in life-science research in both in academia and industry

#### Europe has already paid for the science



Annual cost of generating new protein structure data in labs around the world

Annual cost of maintaining the data in a central database



#### ELIXIR's data infrastructure will...

- Enable full data integration by making the best use of Europe's collective, expanding capacity
- Establish universal principles for optimisation of existing data capacity to meet rising demand - for example, assessing which data should be stored and made available to users
- Present a transparent single interface to a distributed infrastructure



## ELIXIR Data Centre UK funding has been used to purchase:

- - 5-year lease for Tier-3+ space in commercial high-security data centre in London
  - Compute & networking to establish EBI external services in this secure environment
- Benefits of this option:
  - Provides secure service and 24/7 delivery
  - Electricity requirement can be met
  - Scalable
  - Allows just-in-time purchasing
  - Good networking and connectivity



#### ELIXIR's compute infrastructure will...

- Optimise existing compute architecture to enable storage of vast quantities of homogeneous data
- Life scientists produce an enormous quantity of heterogeneous data which needs to be integrated
- An order of magnitude increase in compute infrastructure capacity is needed
- Without it, Europe's data infrastructure will not be adequate to respond to the data deluge



#### ELIXIR's tools will...

- Harmonise the staggering number of analytical software tools
- Support typical analyses using multiple tools linked together into data processing 'pipelines'
- Be developed based on the following principles:
  - discoverability
  - ease of use
  - benchmarking
  - Interoperability



## Maintaining open access Open access to life science data is essential for

- Open access to life science data is essential for advances in many areas of research
- Open access to bioinformatics resources provides a valuable path to discovery, one that in many other areas of research is limited by commercial confidentiality
- Charging for that data, or seeking to restrict access through exercising Intellectual Property (IP) rights, would impede progress
- ELIXIR will guarantee that open access to biological data is maintained
- Speaking with a single voice will strengthen Europe's influence in such global discussions



## **OpenAIRE**



## **ELIXIR bridges life sciences and ICT**

#### Millions of Life Sciences users





#### ICT PROVIDERS

GEANT - Network EGI - GRID PRACE – Supercomputing EUDAT – Data Storage OpenAIRE – OA: Pubs & Data



#### Other ESFRI Research Infrastructures

#### **ELIXIR coordinates BioMedBridges**



## Integrating From molecules to medicine



The BMS Research Infrastructures constitute a well coordinated series of facility providing technologies and expertise relevant in the biological, medical, translational and clinical domains."



#### Building data bridges from biology to medicine in Europe

BioMedBridges is a joint effort of twelve biomedical sciences research infrastructures on the ESFRI roadmap. Together, the project partners develop the shared e-infrastructure—the technical bridges—to allow data integration in the biological, medical, translational and clinical domains and thus strengthen biomedical resources in Europe.







#### **From General to Special**



Network

## **Common ESFRI Reqs for e-RIs**

- Single sign-on: consistent access to resources
- Virtual organisations (collaboration)
- Persistent storage: long-term preservation of data and its access
- Data Management services
- Standards web services
- Workflows support of access to HPC/grid/network resources (compute+data) across Europe
- Training

Global scope: beyond Europe

## **Research Paradigms**

► Hypothesis → Experiment/Observation → Data collection

#### ▶ Data collection → Data analysis → Hypothesis

## **Strategic Decisions Needed**

- Technological
- Architectural
- Financial
- Political
- Legal
- Cultural / Social
- ...al



## **BIG DATA**

- Volume (high)
- Velocity (high)
- Variety (great)
- Veracity (lack of)
- Value (hard to extract)

## Architectural

- Where does data reside?
- Where does computation occur?

- Past/present researcher common behavior
  - discover datasets through metadata queries
  - download all data for local processing
  - keep analysis results locally
- Moving data to computation does not scale
- BIG DATA

## Architectural

#### Future researcher behavior

- discover datasets through metadata and data queries
- upload computation to RI for remote processing
- keep analysis results at RI
- keep computation at RI
- Computation and data e-infras together
- Cultural/social issues
- Political/legal issues



#### Thank you!