

ELIXIR: Data Challenges in the Life Sciences

e-IRG workshop, Athens, 9-10 June 2014

Andrew Smith ELIXIR Hub



ELIXIR's mission

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

bioindustries

environment

society



medicine



The potential



Genome-wide analysis of crop plants

- Population growth and climate change are major challenges to food security.
- Traditional routes to crop improvement are too slow to keep up with this increase in demand.
- Understanding plant genomes helps us identify which species will be most tolerant to drought, salt and pests while still providing optimum nutrition.





Matching the treatment to the cancer

- One in 10 women in the EU-27 will develop breast cancer before the age of 80.
- If we can identify patterns of genes that are active in different tumours, we can diagnose and treat cancers earlier.







The challenges



Growing data



The data challenge: geography

- Data production increasing sites across Europe
- European Illumina sales up 20% 2013



Source: http://omicsmaps.com



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 (commics Databases) Metabolic and Signaling Pathways 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 diatabases Plant diatabases Cell Diological databases

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
Search Summary Papers

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



Users are global



Source: EMBL-EBI Live Data Map



The policy drivers: Open Access to data

- Open access to life science data is essential for advances in many areas of research
- It provides a valuable path to discovery, one that in many other areas of research is limited by commercial confidentiality
- National funders increasingly require researchers to make data open
- EC's H2020 pilot on Open Research Data and Data Management Plans





The response



Infrastructure for Life Sciences





ELIXIR's structure

- Tools
- Standards
- Data
- Compute
- Training
- Industry





ELIXIR Nodes

o elixír ELIXIR: Swedish Node

The Swedish ELIXIR node will initially contribute with the Human Protein Atlas (http://www.proteinatlas.org)



The Swellsh EURIR node will initially contribute with the Human Protein Actis (http://www.perceloritize.org/, which contains: (1) Human Tissue Actis: (2) Human Schollbar Actis; (1) Human Schill Line Actis; and (4) Human Cancer Actis: Facthermone, RNA transcript data has been added for a majority of the tissues in the normal tissue and the cells in the cell line actis. normal tissue attas and the cells in the cell line attas. The Human Procession Adlas (H40) Hydrogramme is ta activitific research programme ied by Port. Mathias Liblier with the goal to explore the whole human persons using an antibody-based gene-centric approach with the effort to map and characterise a representative protein for each protein-coding human gene. New data are released annually to the Human Protein Adlas.

Contact

ELIXIR: The Norway

Node BLS-Bioinformatics

ELIXIR's Norway Node will provide competence and infra structure building on key areas for Norway, in particular marine resources and medical research. Challenges related to processing and analysis of data from next-generation sequencing and other

SolLifeLab - Science for Collaborating organisations

University of Bergen The coordinating partner of the Elsir Norway node. The main focus is on marine genomics and e-influstructure. An early deliverable is LiceBase developed in traffic collaboration withe the Sea Lice Research Centre. University of Oslo Emphasizes biomedical resource provision an analysis, leveraging public resources in integrative statistical genomics, with secure management of person sensitive data.

Norwegian University of Life Sciences Main focus on providing genomic resources for species-oriented an comparative fish genomics. Provision of web-based solutions for se toolboxes, and computational access to these date.

Norwegian University of Science and Technology Tools and resources for analysing genome data, with focus on gene regulation, non-protein-coding RNAs and epigenetics, but also bacte genomics. Handling and analysis of data from furman biobanis.

iversity of Tromse Tools and pipelines for analyzing metagenomic (and genomic) data, with a particular focus on taxonomic classification and bioprospecting (functional and metabolic potential).

Marine research

The Norwegian ELDOR Node will provide services and resources toward marine The Novegian ELDIK Node will provide services and resources toward marine generics is nichtige researchers, povermenter, and industry. The Novegian Node will off several integrated packages general towards large scale analysis of marine genomic and metagenomic data (e.g., fils genomics and marine bein-prospecting). This also includes polylism of web-based solutions for services, toolbases, and computational access to reference data provided by the ELDIK motions.

Health and biobanks

The howay hold supports infrastructure for handling and analysis of data for medical research, including howan holpanks, Such data may be sensible, and must be stored with source access. The nools is developing infrastructure for sensible data. Tools for data analysis are integrated into NACS, the Normegian e infrastructure for life sciences. This provides user friendly solutions for sample for human re-sequencing data and other genome-scale analyses. Ý





Nonvegian Bioinformatics Platform

high-throughput methods are important both to basic researc these areas, and to the development of new enterprises. The node will also provide training and support toward researchers



ELIXIR: Finland Node

The ELROR node in Finland, Biomedinfra.ft, provides compute cloud and strange resources for life sciences with integrated periods are for biomedial organisations and evolved with the de-velopment of the European e infrastructure (EGANY for research resource), EGA and PARCE for computing and EUDAY for data, providers and can be used to host tools and build topical data ser-tements and an be used to host tools and build topical data ser-tements and an a be used to host tools and build topical data ser-anging a science virtualised platform for resear-ies data.

-01

dical data resources are collections and registries opulation that are being organised and digitised Molecular Medicine Finland (FIMM) and the Na-Health and Welfare (THL) and will becme avail National biomedical data interoperation and in-ne reference data of ELDIR provides a use case for vetic diagnostics of, for example, heart diseases,

The ELIXIR node is operated by the Finnish data centre provider CSC-IT Center for Science under Biomedinfra.fi consortium agreen ment with TIMM and THL. It provides most address for a consol-build reference genetic data resolution for training and the formation of the nome sequence variation for training and the formation of t The ICT hardware is hosted by C and university network (Funet) u essary. The node organizes trai The Nordic ELIXIR nodes (Denr



Ministry of Education and Culture ACADEMY OF FINLAND

Biome

Collaborating organisations University o Canine gen Technical University of Denmark University of Copenhagen FinnishCe

Aalto Unive er Science.1 at the Infor

ersity of Southern Denmark

Key services in the Danish ELIXIR node:

ELIXIR: Danish Node

Creation of a comprehensive tool registry

The Danish ELIXIR node on bioinformatics tool inter-

operability and integration will address the main problems in

Creation of a comprehensive tool registry equipped with adequate search functionality based on the existing and emerging ontologies in the area of life sciences. This involves a comprehensive effort in tool description shared by the Node and several ELDOR partners as well as the individual tool providers. The registry will continuously interchange content with other tool registries, existing and emerging.

Provision of tools

el

Legal, ethical and privacy requirements for enabling research on omedical data require solutions for researcher autometication an REMINSTANCE. The Reserves: Enformatic Managemeet Statum (EMM) Second The Reserves Enformation (Managemeet Statum (EMM) Second (Managemeet) (Second Statum) (Seco

Provision 1 of the top petered by user across academia and industry: the Note will encourage and a mint top provident is generated or a faces. Underfaces such gift and by the community, including the constant will also all host general and the changing needs of the cares particular. The Note will active general and the changing needs of the cares particular, the Note will active provide the changing needs of the cares particular, the Note will active to the same same gift here any particular and the same particular to influence assessing the top provident. The Note will active to this and general mining in order to anyone robustions.

Benchmarking, sustainability and renewal

Benchmarking, tool sustainability and renewal: the Node will coordinate the flow of information on the tool quality and relevance, and make it available to the users on regular basis. Automatic update of tools will also be promoted and the Node will actively encourage tool providers to do this.

Tool interoperability

Teol interspeciality: promotion of standard data formats and schemas as well as evaluation and promotion of workflow engines and integrated workbenches. The effort will build upon the origoing efforts in this area within the European covering beind ELDOR.

Contact



Ь

ы

multiple

interfaces

Ashes University

Lundbeck A/S

New North AS



tool utilization experienced by the life sciences community

Glostrup Research Institute, Glostrup Hospital

sive registry precise descriptions searching

2

ATION

DISCOVERY

benchmarking

RELEVANCE

evaluation

Novorumes A/S

Exigon A/S



bio registry



Partners About Home

Genomics	Q			Show column	s [*] Clear filter
Name 🔺 🔻	Туре	Description	Topics 🗸	Web UI 👻	REST API
	Filters 👻				
ANNOVAR	Tool	variant annotation	Genomics		
Arabidopsis nucleolar protein databa	Database	Comparative analysis of human and Arabidopsis nucleolar proteomes	Genome, proteome and		
ArrayExpress	Database	A database of functional genomics experiments including gene expression w	Functional genomics	\checkmark	
ButterflyBase database of comparativ	⊻ Database	Gene database for butterflies and moths.	Invertebrates Comparati		
CDinFusion	Tool	A submission-preparation-tool for the integration of contextual data (CD) wi			
Comprehensive microbial genome re-	Database	Information on all of the publicly available, complete prokaryotic genomes	Prokaryotes and archae		
Comprehensive yeast genome databa	Database	Information on the molecular structure and functional network of the entirel	Pathways, networks and		
CREST	Tool	CREST (Classification Resources for Environmnetal Sequence Tags) is a colle	Metagenetics Amplicon		
Database of homologous sequences	Database	Homologous genes from fully sequenced organisms. It allows selection of se	Nucleic acid sequence al		
EBI metagenomics	Database	Pipeline for analysis of metagenome sequences	Metagenomics		
Ensembl Genomes Bacteria	Database	Ensembl Genomes bacteria division	Genomics		
Ensembl Genomes Fungi	Database	Ensembl Genomes fungi division	Genomics		
Ensembl Genomes Plants	Database	Ensembl Genomes plants division	Genomics		
Ensembl Genomes Protists	Database	Ensembl Genomes protists division	Genomics		
Ensembl Genomes	Database	Extension of Ensembl genome database across the taxonomic space	Genomics		
Ensembl Human	Database	Ensembl human	Genomics		
Ensembl Metazoa	Database	Ensembl Genomes metazoa division	Genomics		
Ensembl Mouse	Database	Ensembl mouse division	Genomics		
Ensembl REST API service	Service endpoint	Retrieves Gene Tree dumps for a given Gene Tree stable identifier	Comparative genomics		\checkmark
Ensembl REST API service	Service endpoint	Retrieves the Gene Tree that contains the given stable identifier	Comparative genomics		\checkmark
Ensembl REST API service	Service endpoint	Retrieves a Gene Tree containing the Gene identified by the given symbol	Comparative genomics		\checkmark
Ensembl REST API service	Service endpoint	Retrieves homology information by ensembl gene id	Comparative genomics		\checkmark
Ensembl REST API service	Service endpoint	Retrieves homology information by symbol	Comparative genomics		\checkmark
Ensembl REST API service	Service endpoint	Lists all available comparative genomics databases and their data release	Information		\checkmark
Ensembl REST API	Service	Ensembl REST API endpoints	Genomics		\checkmark

BioMedBridges Software Tools Registry



►





ELIXIR pilots addressed key challenges in biomedical research

- Cloud computing "Embassy cloud": Access reference data in a virtual environment – work as though you are at EMBL-EBI or SIB, Switzerland
- Authentication & Authorisation Improved methods and processes for access to clinical data





Identifying new drug targets

ELIXIR pilot: Interoperability of high-resolution protein data at EMBL-EBI and HPA, Sweden





The Human Protein Atlas portal is a publicly available database with millions of high-resolution images showing the spatial distribution of proteins in 46 different normal human tissues and 20 different cancer types, as well as 47 different human cell lines.



European ELIXIR Data - "LightPath" (EBI / CSC)

- To explore the replication of large scale (Petabyte scale) archives to remote sites
- To create a separate source of data files for challenging DataIO projects
- Selection of pilot data transfer technology between EBI and CSC
- Established a dedicated light path between datacenters in London and Kajaani
- Development of model for future IO needs in the life sciences in Europe





Cross-site VM Operation - pilot

- Perform analysis via cloud infrastructures and VMs
- Transfer VMs between computing centers to allow researchers to perform analyses that they could not otherwise do locally
- Supported by 5 NRENs and in collaboration with



ENLIGHTEN YOUR RESEARCH GLOBAL	janet FUET	
Home About Projects How to submit Important dates Partners News	₽ Search	
Contact		
International Networks to Aid Global Research Collaborations in Climate,	MAILING LIST Do you want to receive an email notification for the next Enlighten Your Research Global Call for Proposals? Please sign up for the	
Bioinformatics and Computer Science Posted on 19 November 2013 by Mary Hester	mailing list!	
INTERNATIONAL "ENLIGHTEN YOUR RESEARCH" PROGRAM SELECTS FOUR GLOBAL DATA-INTENSIVE RESEARCH PROPOSALS	Last name Email	



Cross-site VM Operation



European Research Infrastructures





Knowledge exchange workshop

Discussion of big data challenges in life sciences

- Focus on few representative domains
- Looking 5 years ahead
- Jointly identify potential solutions to our problems







Thank you

