# ELIXIR

## Technical activities in ELIXIR Europe
### *Rafael C Jimenez*
*ELIXIR CTO*

*EGI forum 2014, Tuesday 20 May*

*European Life Sciences Infrastructure for Biological Information*
*www.elixir-europe.org*

# China High Speed Rail Network

# European Research Infrastructures

# European Research Infrastructures

# Research infrastructures

## Facilitate research

Physical facilities
Scientific information

**LS** **ICT**
life sciences e-infrastructures

Transfer
Computation
Storage

Data

# Data resources in life science

- **Many**
- **Diverse**
- **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%)

# Utility of databases



Scientific impact

Too little
information

Too many databases
Too diverse interfaces

# Data interoperability

# Improving Links Between distributed European resources

## ELIXIR pilot: Interoperability of protein expressions resources
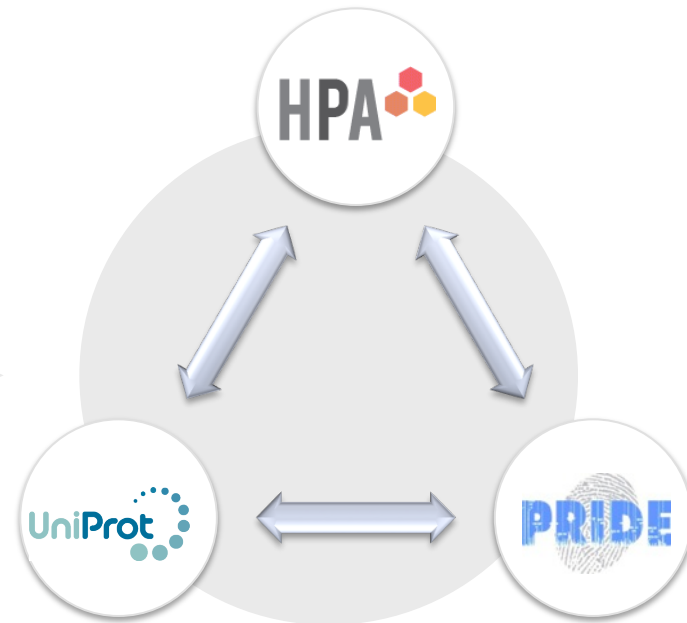


The Human Protein Atlas portal is a publicly available database with millions of high-resolution images showing the spatial distribution of proteins in 46 different normal human tissues and 20 different cancer types, as well as 47 different human cell lines.

# Growing data

# Proteomics data in PRIDE

# ~85% raw data

# Data types examples

| Raw data | Process data | Metadata |
|---|---|---|
|  | TTGTTATCCG… | *DNA*<br>*Human*<br>*Liver*<br>*Mitochondria*<br>*W. Smith*<br>*…* |
|  | LPISASHSSK… | *Peptide*<br>*Mouse*<br>*Heart*<br>*Nucleus*<br>*J. Heinz*<br>*…* |
| *…* | *…* | *…* |

# Data submission

raw data

processed data

metadata

Submissions

Centralized database

# Data submission - pilot

EUDAT

raw data

PID

Submissions

processed data

metadata

# Data analysis

# 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

# Cross-site VM Operation

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

- Aim
  - 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
- Update:
  - 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 lifesciences in Europe

GÉANT

CONNECT

THE MAGAZINE FROM THE GÉANT COMMUNITY | ISSUE 13 OCTOBER 2013
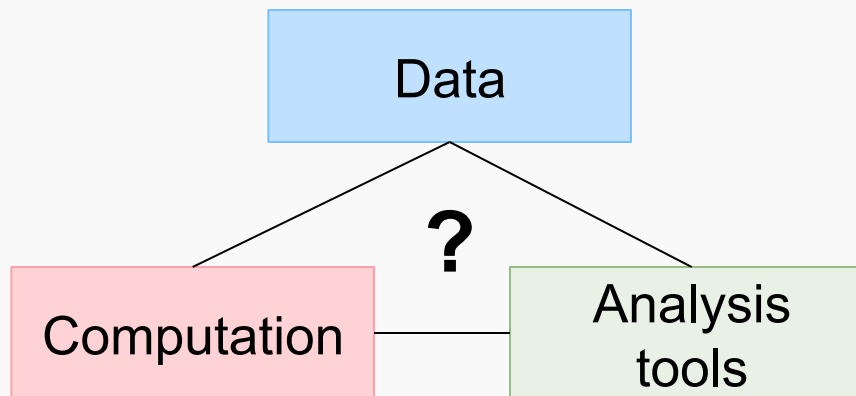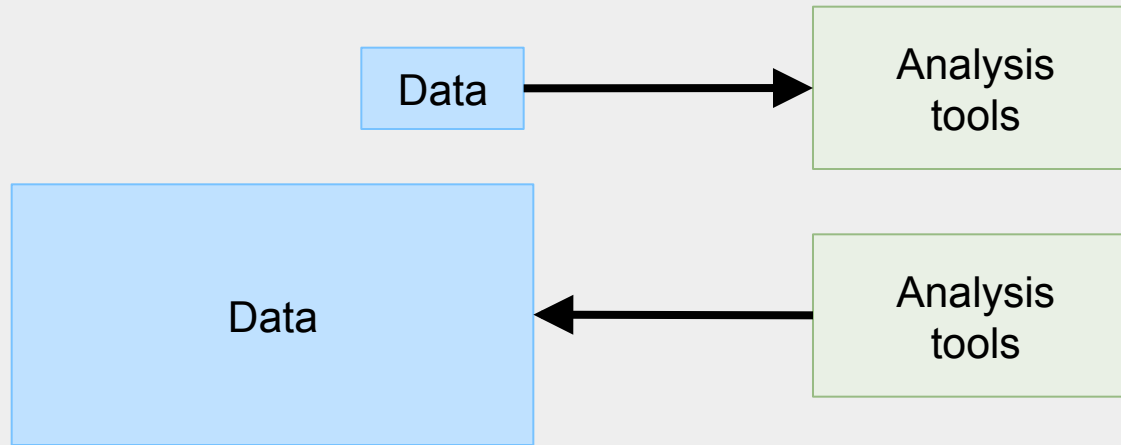
DNA
DATA
DELUGE

GÉANT AND
THE GLOBAL
SHARING OF
GENOMIC
DATA

CREATING NETWORKS:
BEHIND THE SCENES AT
GÉANT'S NETWORK
MIGRATION

BUILDING COMMUNITIES:
VIA NETWORK
PERFORMANCE AND
MONITORING

SUPPORTING NRENS:
CELEBRATING 20 YEARS OF
EENET, THE INTERNET AND
INDEPENDENCE

# REMS - Resource Entitlement Management System

- Access to sensitive data (genomics) granted by a Data Access Committee

- In collaboration with eduGAIN

- Agreements to be applied to other domains: FI-CLARIN & FI-CESSDA

Starting point: https://remsdemo.csc.fi/

# The Economist

World politics | Business & finance | Economics | Science & technology | Culture | Blogs | Debate | Multimedia | Print edition

**Business**

# Welcome to the yotta world

Comment (1)    Print

E-mail    Reprints & permissions

**Big Data will flood the planet**

Nov 17th 2011 | From The World In 2012 print edition

Like 205    Tweet 206

**Exaponential**
Quantity of global digital data, exabytes

1,000 (kilo)

1,000,000 (mega)

1,227
2012

2,720

7,910
2015

2005

Source: EMC/IDC Digital Universe Study, 2011

Even if you still have to think twice about the meaning of "giga" and "tera" in computer-speak, you'd better get ready for "peta", "exa" and "zetta". These binary prefixes, which

For Big Data to become huge, however, there are still hurdles to leap. For one thing, the tools to analyse data are not yet good enough. And **people with the skills to analyse data are scarce and will become scarcer**. By 2018 there will be a "talent gap" of between 140,000 and 190,000 people, …

Follow *The Economist*

Latest blog posts - All times are GMT

**The Economist explains**: How America defines religious freedom

**Poland's agriculture**: A golden age for Polish farmers?

**The enigma of flight 370**: Dashed hopes
Newsbook - Mar 24th, 17:56

Feedback

elixir

20

# Challenges

- **Sustain** data and services
- Make data and service **interoperable**
  - Necessary to integrate data
  - Specially medical, clinical and research
- Data too big to **store**, **exchange** & **compute**? Forthcoming challenges …
  - Data production grows faster than storage
  - Cost of data production technologies declines faster than storage
  - It takes longer to transfer data than produce the data.
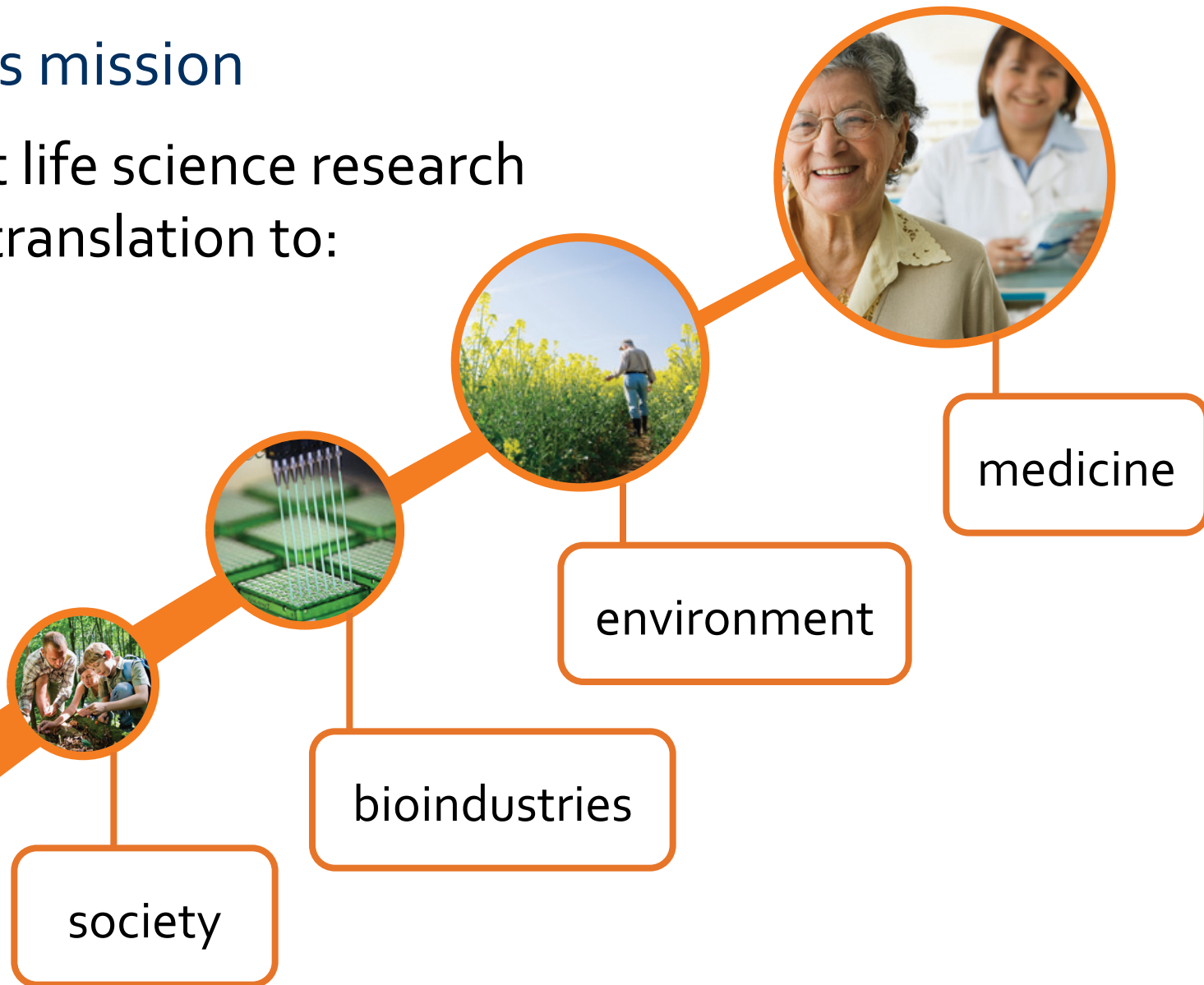- Privacy, security & access (AAI)
- Training

# ELIXIR

- European **life sciences** research infrastructure for **biological information** to **facilitate research**

- **Safeguard data** and build **sustainable data services**

- Participated by major bioinformatics service providers and supported by **17 EU member states**

- Creating a robust infrastructure for biological information is a **bigger task than any individual organisation** or nation can take on alone
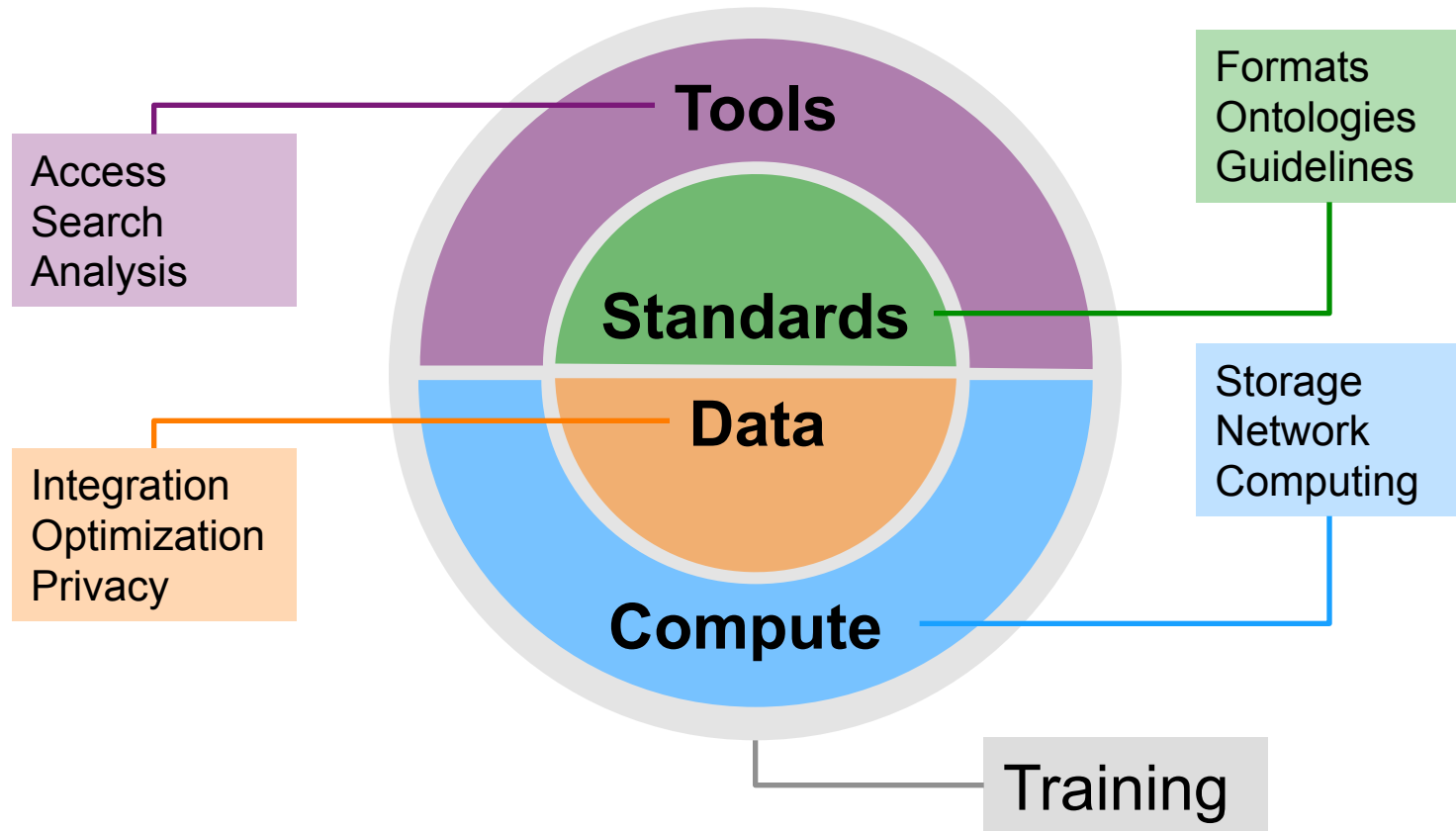
# ELIXIR's mission

Support life science research and its translation to:

medicine

environment

bioindustries

society

# Infrastructure for Life Sciences



Tools

Standards

Data

Compute

Access
Search
Analysis

Formats
Ontologies
Guidelines

Integration
Optimization
Privacy

Storage
Network
Computing

Training

# Acceleration towards sustained operations

| | | | PERMANENT PHASE | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 |

**PREPARATORY PHASE** — **INTERIM PHASE** — **CO-ORDINATION** — **SERVICE DEPLOYMENT** — **SUSTAINED OPERATIONS**
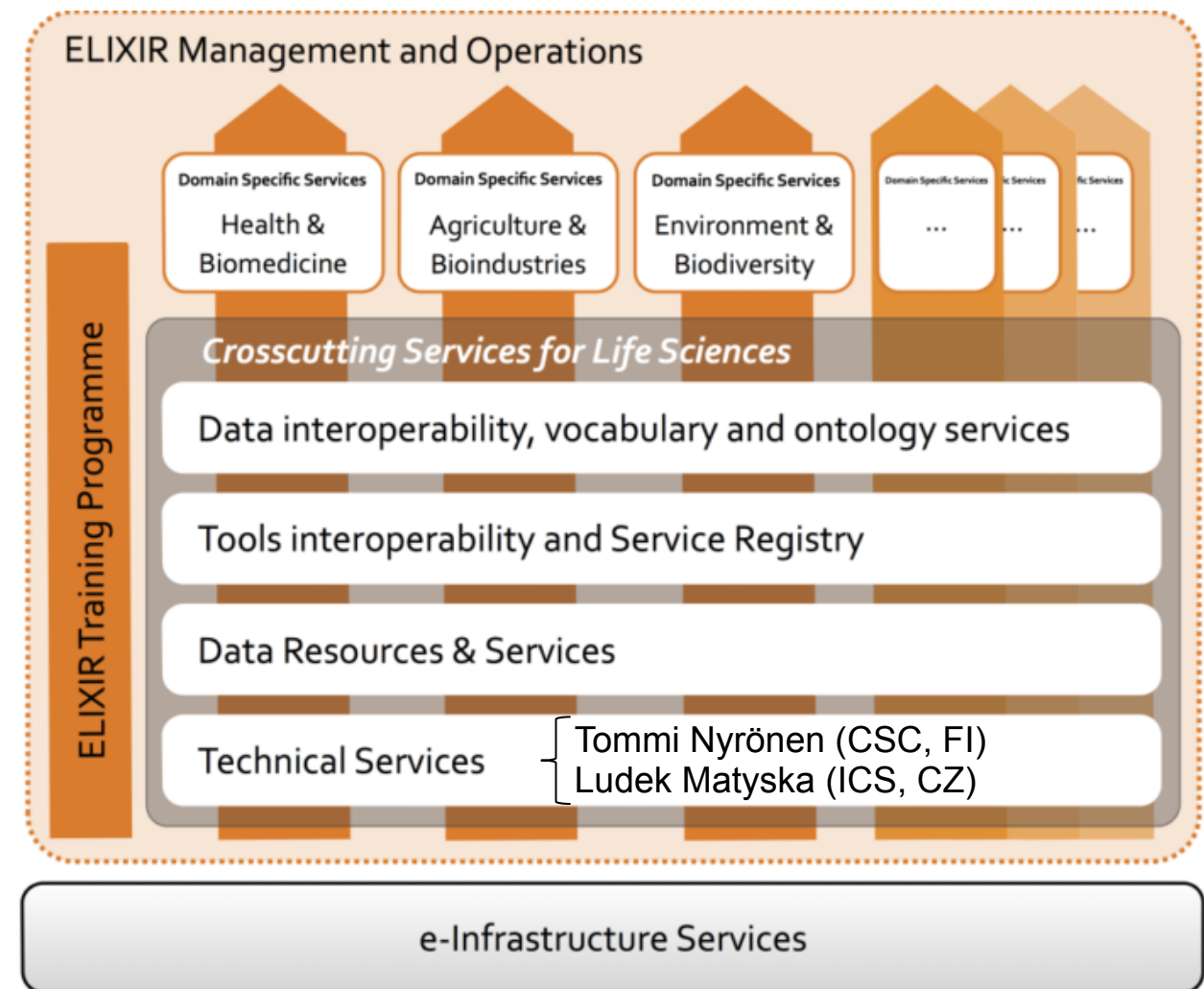
- Key objectives:
  - In 2014: Build on success of Interim Phase, formal establishment of Nodes
  - In 2015: Start delivering ELIXIR services
  - From 2016 onwards: sustained operations across ELIXIR

# Work streams

# TCG - Technical Coordinators Group

- Coordinate the ELIXIR technical strategy
- Composed of technical representatives from each node

| | | | |
|---|---|---|---|
| **Belgium** | Lieven Sterck | **Norway** | Kjell Petersen |
| **Czech Republic** | David Antos and Jan Paces | **Portugal** | TBC |
| **Denmark** | Kristoffer Rapacki | **Slovenia** | Brane Leskosek |
| **Estonia** | Hedi Peterson | **Spain** | Victor de la Torre |
| **Finland** | Jarno Laitinen and Olli Tourunen | **Sweden** | Mikael Borg |
| **France** | Christophe Blanchet | **Switzerland** | Heinz Stockinger |
| **Greece** | TBC | **UK** | Mario Caccamo and Manuel Corpas |
| **Israel** | Jaime Prilusky | **EMBL** | Steven Newhouse |
| **Italy** | Federico Zambelli | **ELIXIR Hub** | Rafael C Jimenez |
| **Netherlands** | Rob Hooft | | |

# Task forces

- Plan, agree and implement technical strategies
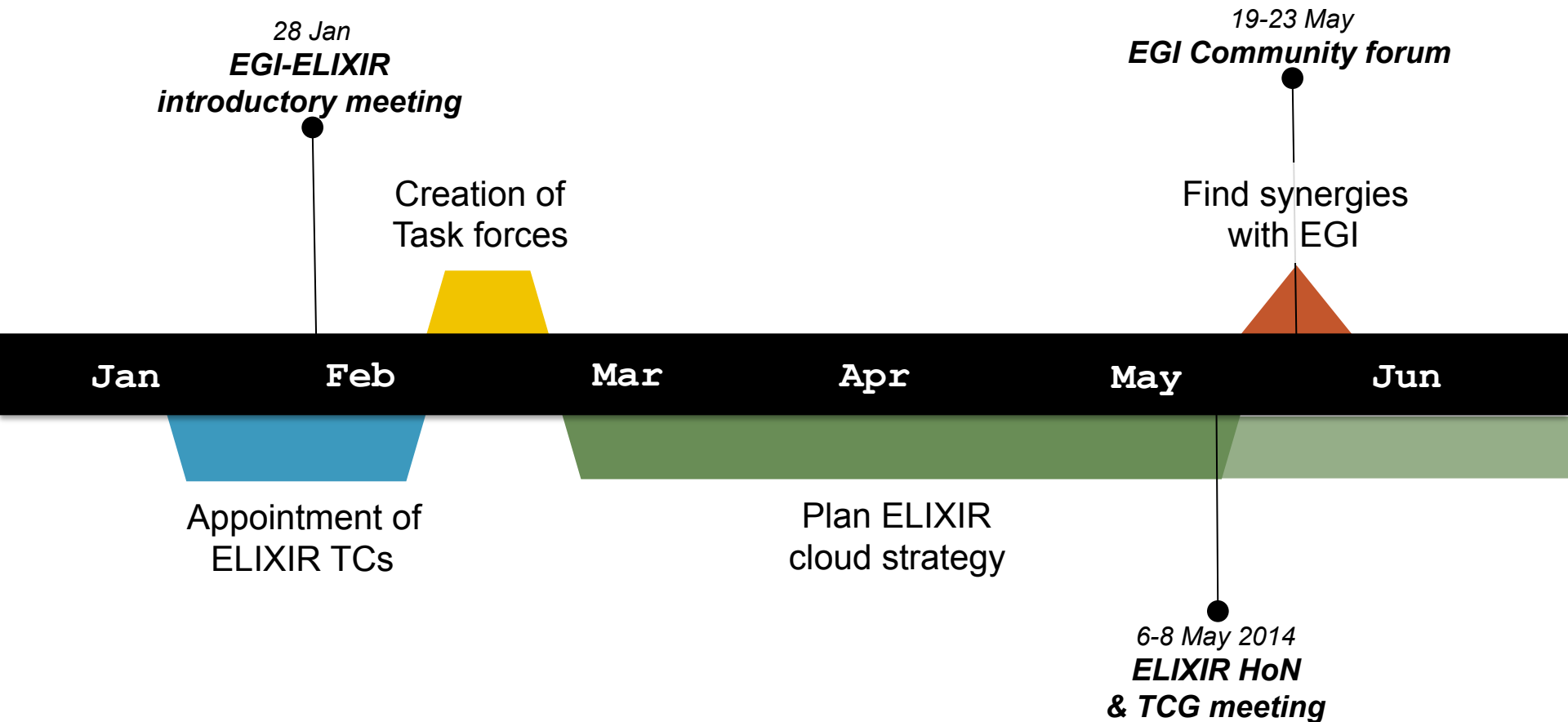- Represent ELIXIR on one specific technical topic

# Prioritized task forces

- **Cloud**
- Service registry
- Communication
- Metrics, monitoring & quality control
- Website
- **Storage**
- **Authentication and authorisation**
- Training portal
- e-Learning and Training

# ELIXIR Cloud task force timeline

*28 Jan*
**EGI-ELIXIR introductory meeting**

*19-23 May*
**EGI Community forum**

Creation of Task forces

Find synergies with EGI

| Jan | Feb | Mar | Apr | May | Jun |
|-----|-----|-----|-----|-----|-----|

Appointment of ELIXIR TCs

Plan ELIXIR cloud strategy

*6-8 May 2014*
**ELIXIR HoN & TCG meeting**

# Thank you for your attention

*European Life Sciences Infrastructure for Biological Information*

*www.elixir-europe.org*

# ELIXIR Pilot Projects

- Five short-term Pilot Actions are underway to act as **test beds for integration of ELIXIR services**:

1. ELIXIR Facing **Cloud Support and Virtual Machines** - with SIB

2. ELIXIR Data IO to pilot the **continuous transfer of major archive resources** to a remote European location - with CSC, Finland

3. Establishing EGA **Distributed authentication** - with CSC, Finland

4. Establishing **EGA** as joint venture – with CRG, Spain

5. **Improving links** between Human Proteome Atlas (HPA) and EMBL-EBI resources

# European Genome-phenome archive

EGA data growth

- EGA is to be **distributed effort** with **archive**, **submission**, and **data distribution** capacity at both the EBI and CRG

- From the users point of view, EGA remains one **integrated** Archive of **secure human biomedical research data**.

- **Search** of datasets at either website is "**global**" across the EGA.

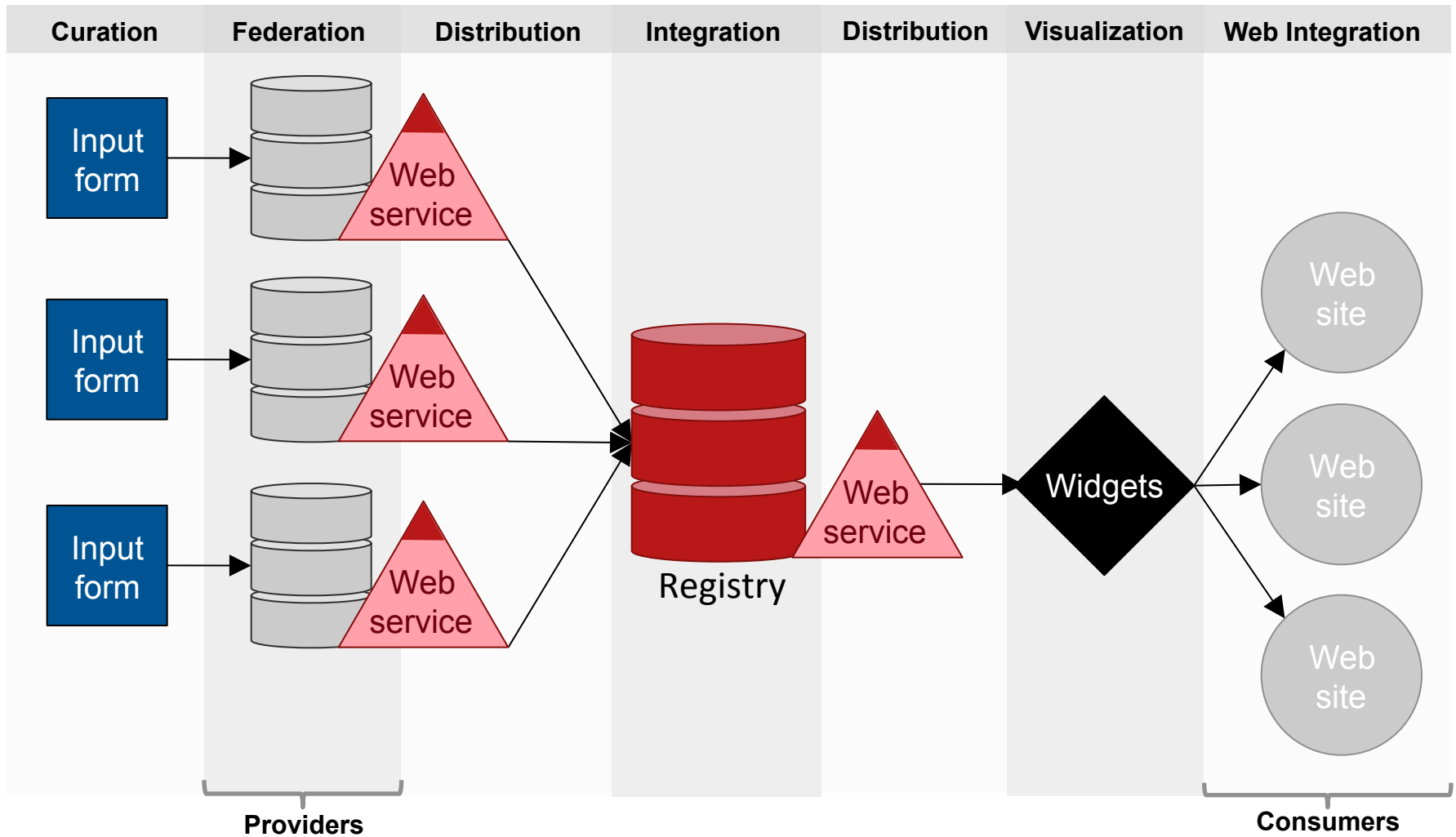# CASE: process for applying access to the Nordic Control Database

# ELIXIR pilots to address key challenges in biomedical research:

1. **Cloud computing**
   "**Embassy cloud**": Access reference data in a virtual environment – work as though you are at EMBL-EBI or SIB, Switzerland

2. **Authentication & Authorisation**
   Improved methods and processes for access to clinical data

3. **High-Performance Computing**
   "**Lightpath**": Connections for on-demand reference data to remote HPC centres at EMBL-EBI and CSC Finland

# Registry of Services and databases

| Curation | Federation | Distribution | Integration | Distribution | Visualization | Web Integration |
|---|---|---|---|---|---|---|

Input form

Web service

Input form

Web service

Input form

Web service

Registry

Web service

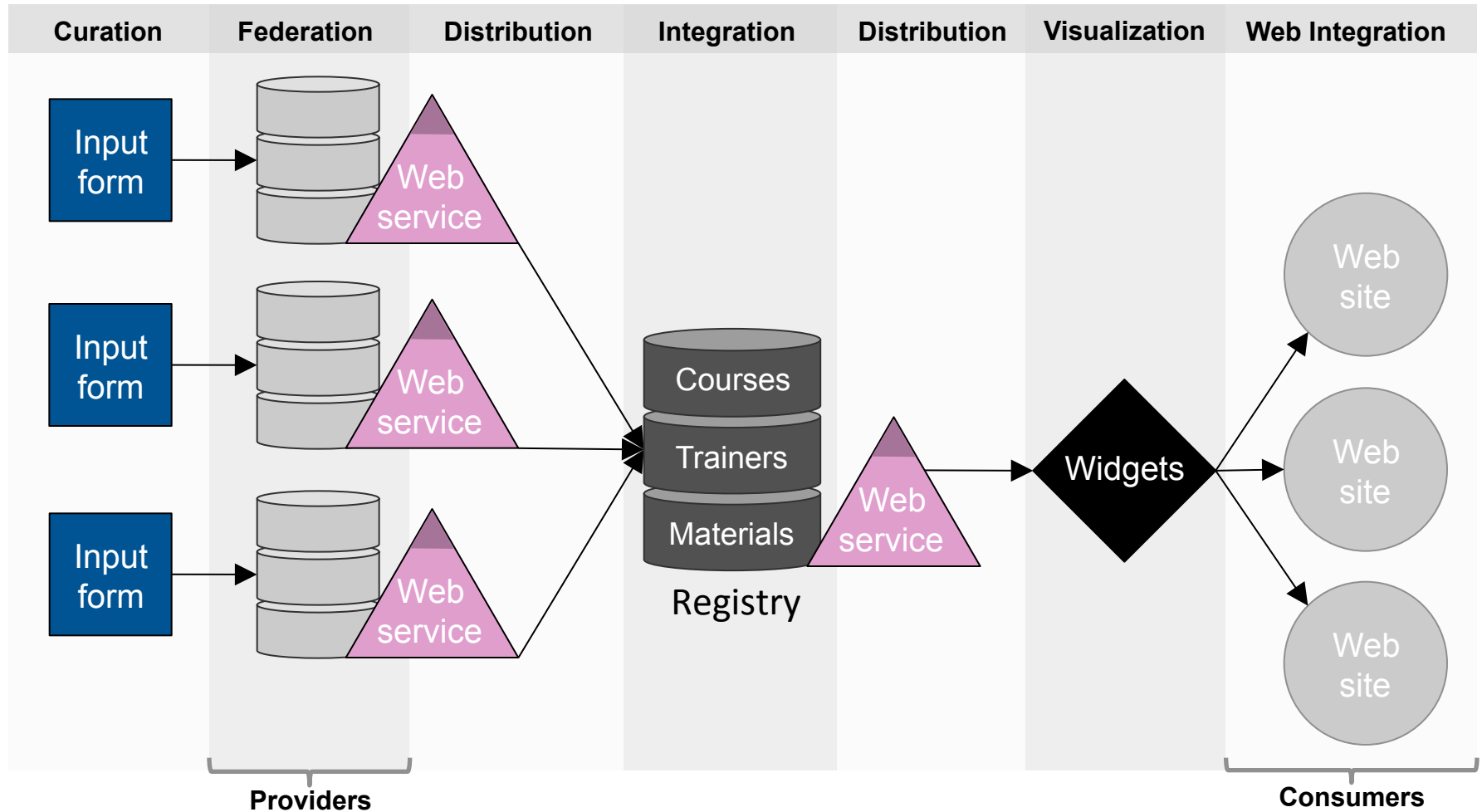Widgets

Web site

Web site

Web site

**Providers**

**Consumers**

Form used by providers to collect data

Data provider

Centralized repository integrating data

Common query interface (web service) based on an agreed standard

Widgets to present information

Websites interested to integrate information from the centralized

# TeSS
## Registry of courses, trainers and materials

# Minimum metadata

life sciences

- Title
- Description
- Creator
- Publication Date
- **Topics**
- **Audience**
- …

Courses

Materials

Trainers

…

Services

Jobs