# ELIXIR: data for molecular biology and points of entry for marine scientists

*Guy Cochrane, EMBL-EBI*

EuroMarine 2018 General  Assembly meeting

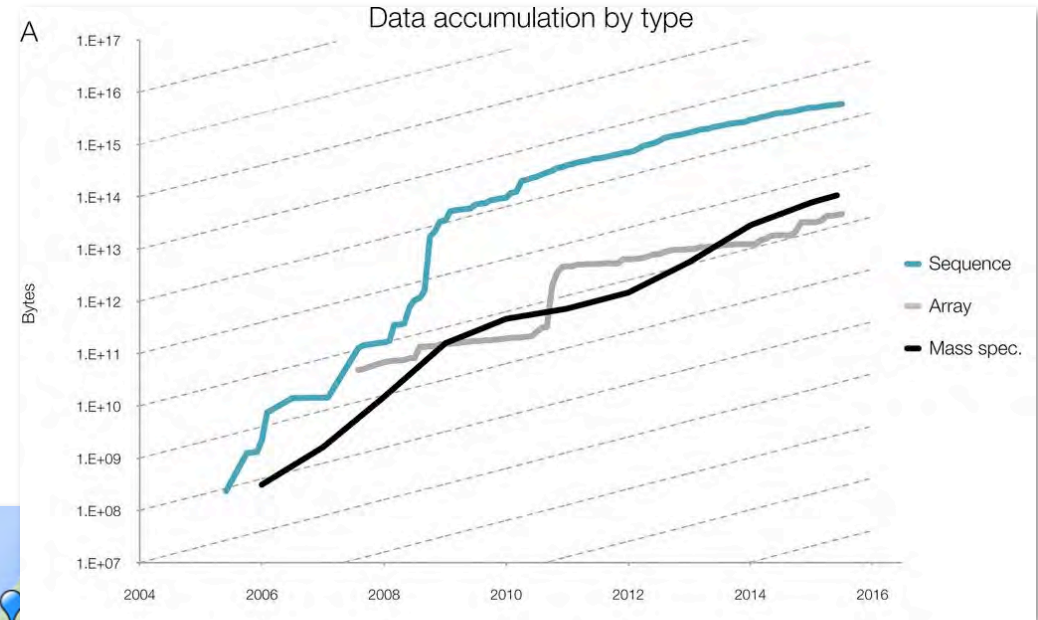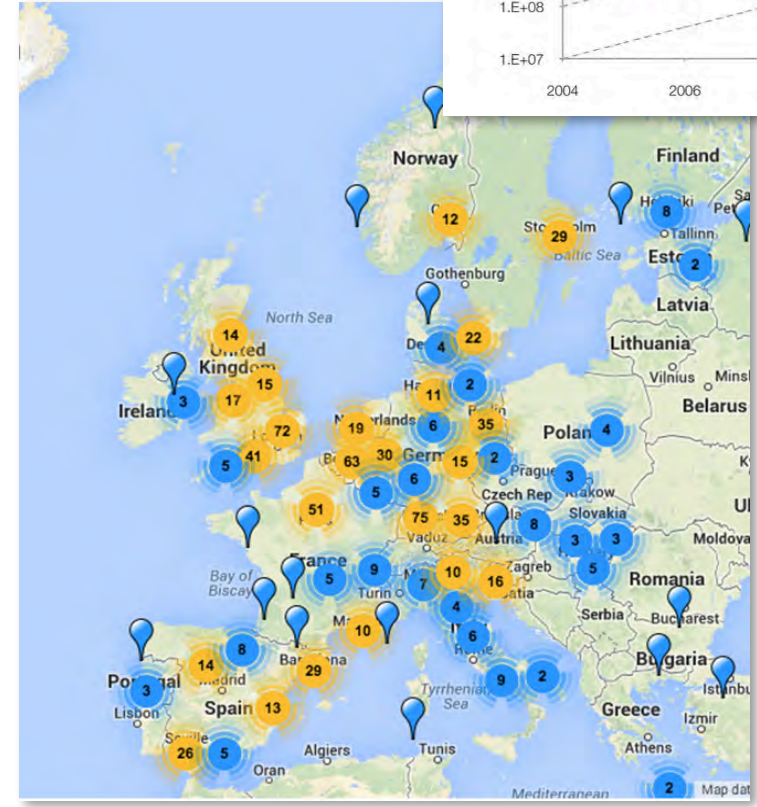17-18 January 2018

*www.elixir-europe.org*

# Scales of molecular biology data



Credits:
Applied Biosystems
Oxford Nanopore Technologies
https://universe-review.ca/R11-16-DNAsequencing.htm


Data accumulation by type

Source: Charles E. Cook et al. *Nucl. Acids Res.* 2016;44:D20-D26




*Source: http://omicsmaps.com*

# Data resources in life science



Non-vertebrate genomics

Structure

Metabolic and signalling p...

Nucleotide sequence

Protein sequence

RNA sequence

...genom...

Other

Human genes and diseases

Plant

Microarray and...

Organelle

Immunol...

Proteomics

## 1800
- Diverse
- Plentiful
- Dispersed

**databases**

elixir

# ELIXIR Members

Belgium

Czech Republic

Denmark

EMBL

Estonia

Finland

France

Germany

Hungary

Ireland

Italy

Israel

Luxembourg

Netherlands

Norway

Portugal

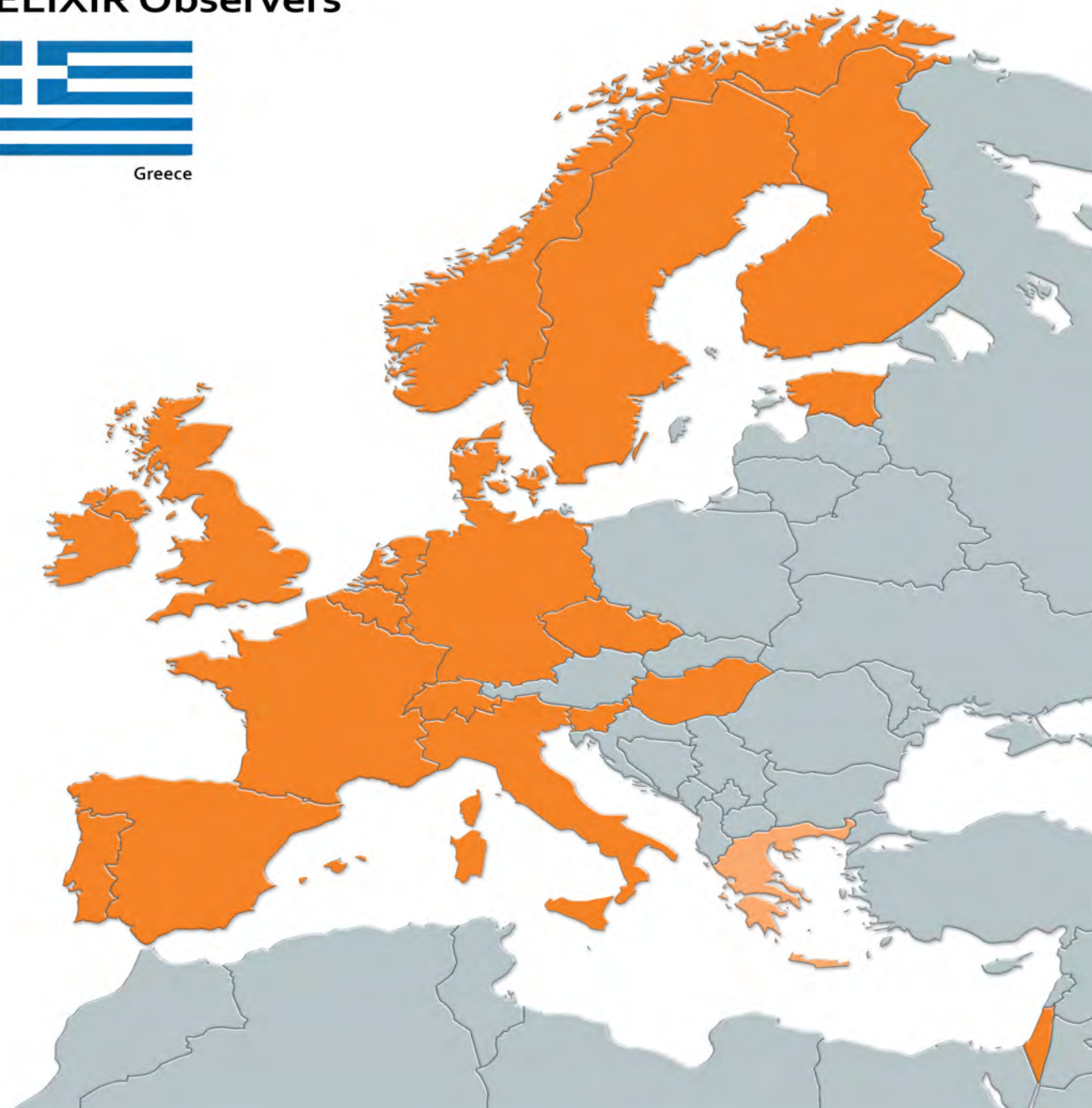Slovenia

Spain

Sweden

Switzerland

United Kingdom

# ELIXIR Observers

Greece

# ELIXIR: European infrastructure for biological information

**Data infrastructure** for Europe's life-science research:
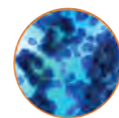
- Data
- Interoperability
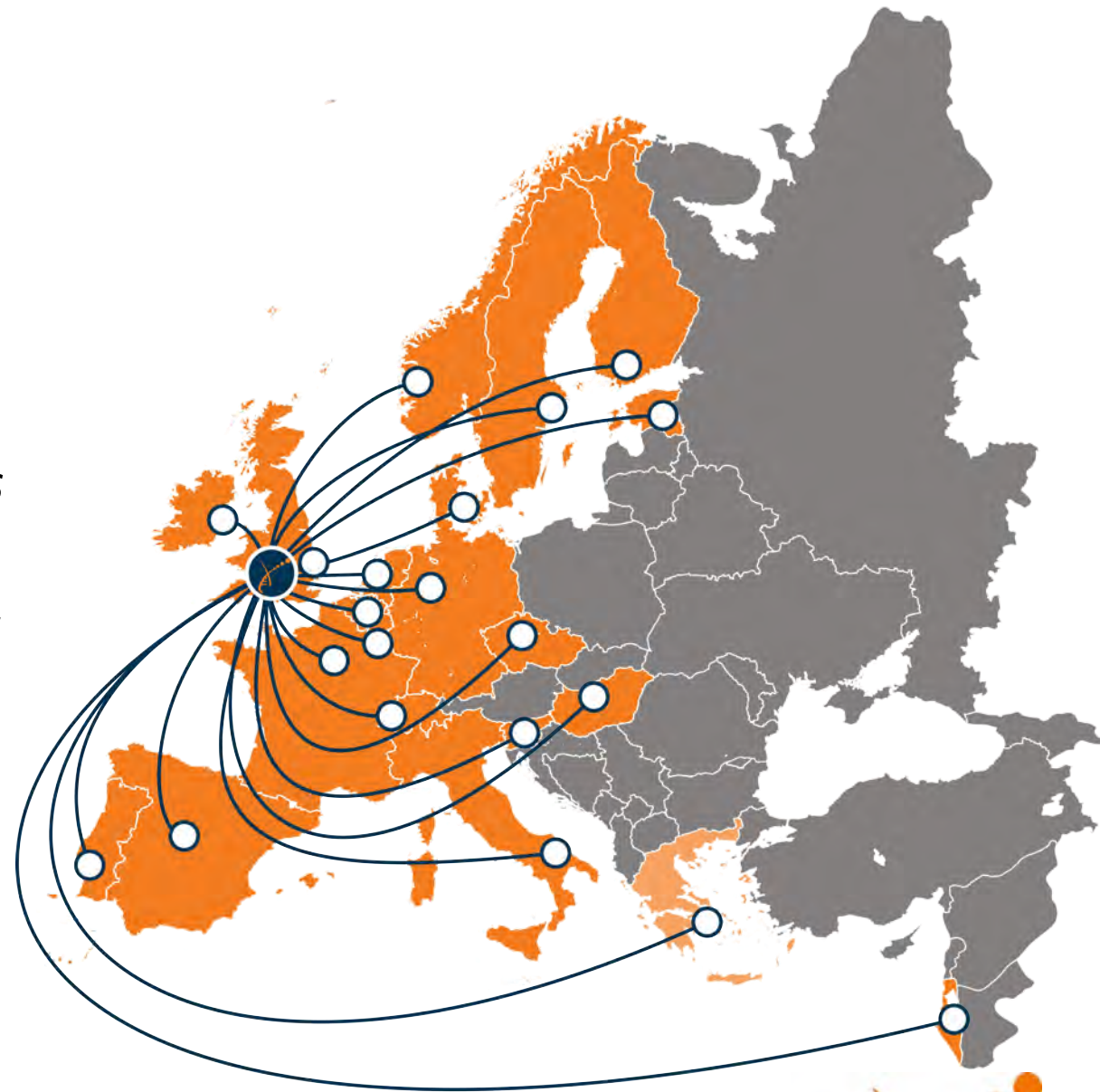- Tools
- Compute
- Training

 Marine metagenomics

 Crop and forest plants

 Human data

 Rare diseases

@ELIXIREurope     www.elixir-europe.org

elixir

# ELIXIR Services

**Data deposition:**
ENA, EGA, PDBe, EuropePMC, …

**Data management:**
Genome annotation
Data management plans

**Added value data:**
UniProt, Ensembl, OrphaNet, …

**Data Interoperability:**
BioSharing, identifiers.org and OLS

**Compute:**
Secure data transfer, cloud computing, AAI

**Bioinformatics tools:**
Bio.tools

**Industry:**
Innovation and SME programme
Bespoke collaborations

**Training:**
TeSS, Data Carpentry, eLearning

elixir

# Core Data Resources – a data infrastructure

- Scientific focus and quality of science
  - Curation level, benchmarking

- Community served by the resource
  - Web statistics

- Quality of service
  - Uptime, user support and training

- Legal and funding infrastructure
  - Institutional support, use policy

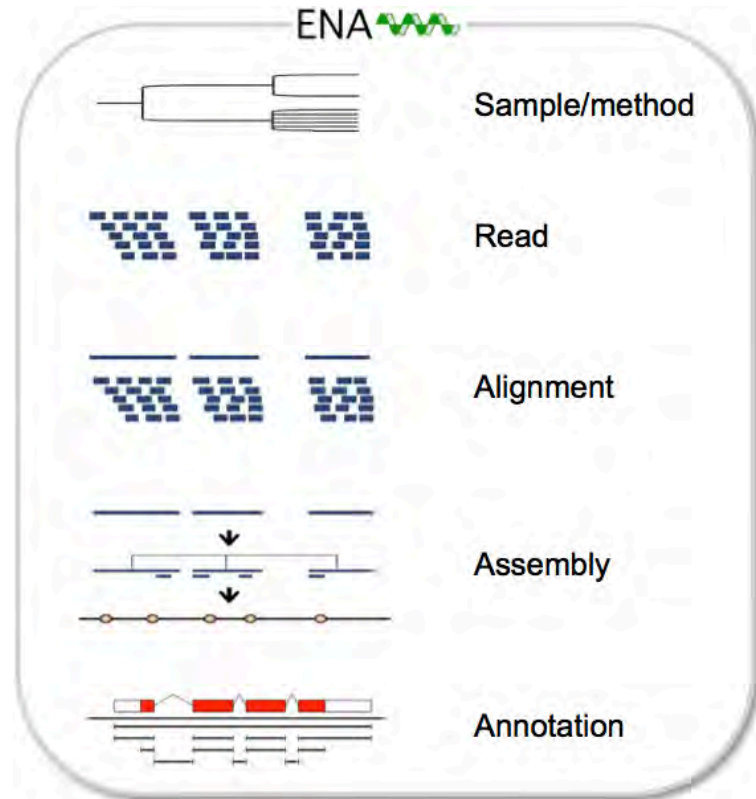- Impact and translational stories
  - Foundational role

Initial set of Core Data Resources

- ArrayExpress
- ChEBI
- ChEMBL
- EGA
- ENA
- Ensembl
- Ensembl Genomes
- Europe PMC

- Human Protein Atlas
- The IMEx Consortium (IntAct and MINT)
- InterPro
- PDBe
- PRIDE
- STRING db
- UniProt

Durinx C, McEntyre J, Appel R *et al.* Identifying ELIXIR Core Data Resources *F1000Research* 2016, **5**(ELIXIR):2422

# Core resource example: European Nucleotide Archive (ENA)



ENA

http://www.ebi.ac.uk/ena/

- **Globally comprehensive scientific record** and European node of INSDC

- A **broad platform** for the management, sharing, integration and dissemination of sequence data

- Established in the early 1980s, extended for **new technologies and applications**

- **Connectivity** with broader EMBL-EBI resources

- Sequence data **foundation**

- **Sustained** within EMBL-EBI under EMBL funding with additional support from EC, UK Research councils, Wellcome Trust, etc.

- **Substantial scale**: 1 submission every 6 minutes, 1.3 petabase pairs across 1.5 million taxa, 2,000-5,000 active data providers, global consumer userbase

- Rich submission, discovery and retrieval **software, tools and services**

EMBL-EBI    elixir

# Use Case: Marine metagenomics

- Mission: develop sustainable metagenomics infrastructure to enhance research and innovation

- Actions:

  - Develop standards for the marine domain

  - Develop databases specific to marine metagenomics

  - Implement tools and pipelines for metagenomics analysis

  - Develop search engine for interrogation of marine metagenomics datasets

  - Develop and deliver training for end users

https://www.elixir-europe.org/use-cases/marine-metagenomics

**Nils P. Willassen ,** NO & **Rob Finn**, EBI

# Tara Oceans expeditions: facts and figures

## 35,000 samples

Scientists on the expedition helped collect, pack and ship around 35,000 samples of plankton and water.

## 7012 datasets

One of the richest molecular sample collections in the public domain - freely available to everyone.

## 11,535 gigabytes

Size of the *Tara* datasets in the European Nucleotide Archive as of May 2015. This represents 12,581 gigabases - roughly equivalent to 135 fully sequenced human genomes.

## Unlimited

Potential to discover new things about life in the world's oceans.

## 40 million genes

Largest-ever DNA sequencing effort for ocean science.

Genetic sequences collected could represent tens of thousands of new species and ecosystem interactions.

Considering the size of the world's oceans, there is much, much more to discover.

*Tara* data: www.ebi.ac.uk/services/tara-oceans-data

EMBL-EBI

# EMBRIC



- **European Marine Biological Research Infrastructure Cluster to promote the Blue Bioeconomy**

  - Objectives

    - Coherent chains of high quality services across RIs

    - Strengthen connection between science and industry

    - Defragment communities

    - TNA - 15 February 2018 opening of the 2nd call (http://www.embric.eu/access/TA)

  - Consortium

    - Led by UPMC – Bernard Kloareg

    - 27 partners

    - INFRADEV-4 'RI cluster' call

# EMBRIC Configurator



- 22 registered consultants

- 2 ongoing cases

- 4 complete cases

- Call for:
  - More consultants
  - More cases

# EMBRIC Configurator examples

- **Sustainability**: Long-term operation of domain-specific data resource on fish species

- **Data dissemination**: *Pseudo-nitzschia multistriata* genomics, transcriptomics and gene expression

- **Data access**: Small molecule activities with respect to strains and genes

# Acknowledgements

# Talk to us

Annalisa Milano
Data Coordination
Biocurator

Jeena Rajan
Sequence Data
Biocurator

Guy Cochrane
Team Leader