ELIXIR Cloud and Compute and the
Global Alliance for Genomics and Health (GA4GH)
ICRI EGI: Enabling global research with interoperable digital infrastructures
19 October 2022

**Jonathan Tedds (ELIXIR Hub)**
ELIXIR Compute, Tools Platform & EOSC Coordinator

jonathan.tedds@elixir-europe.org        *www.elixir-europe.org*

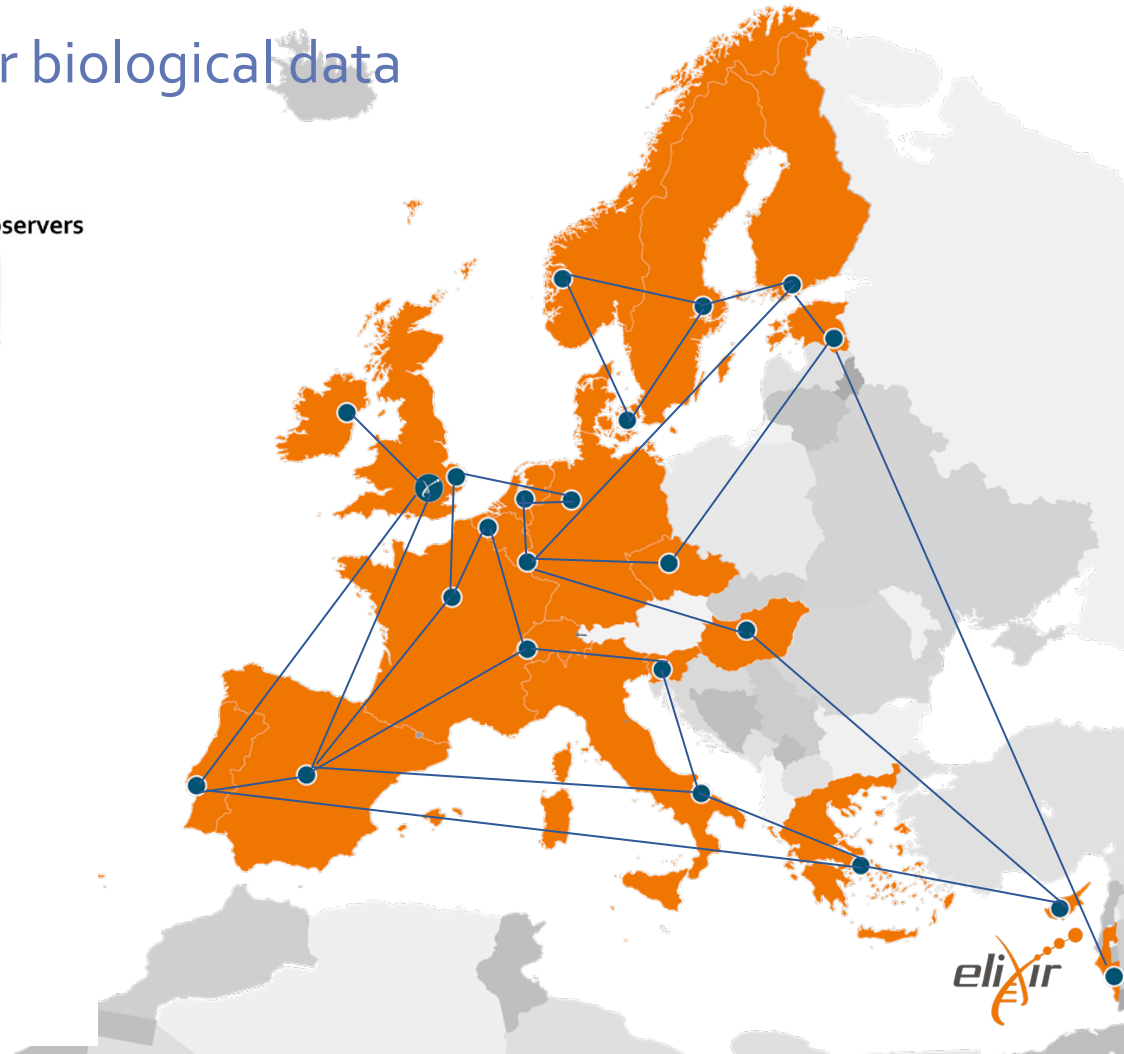# A sustainable infrastructure for biological data
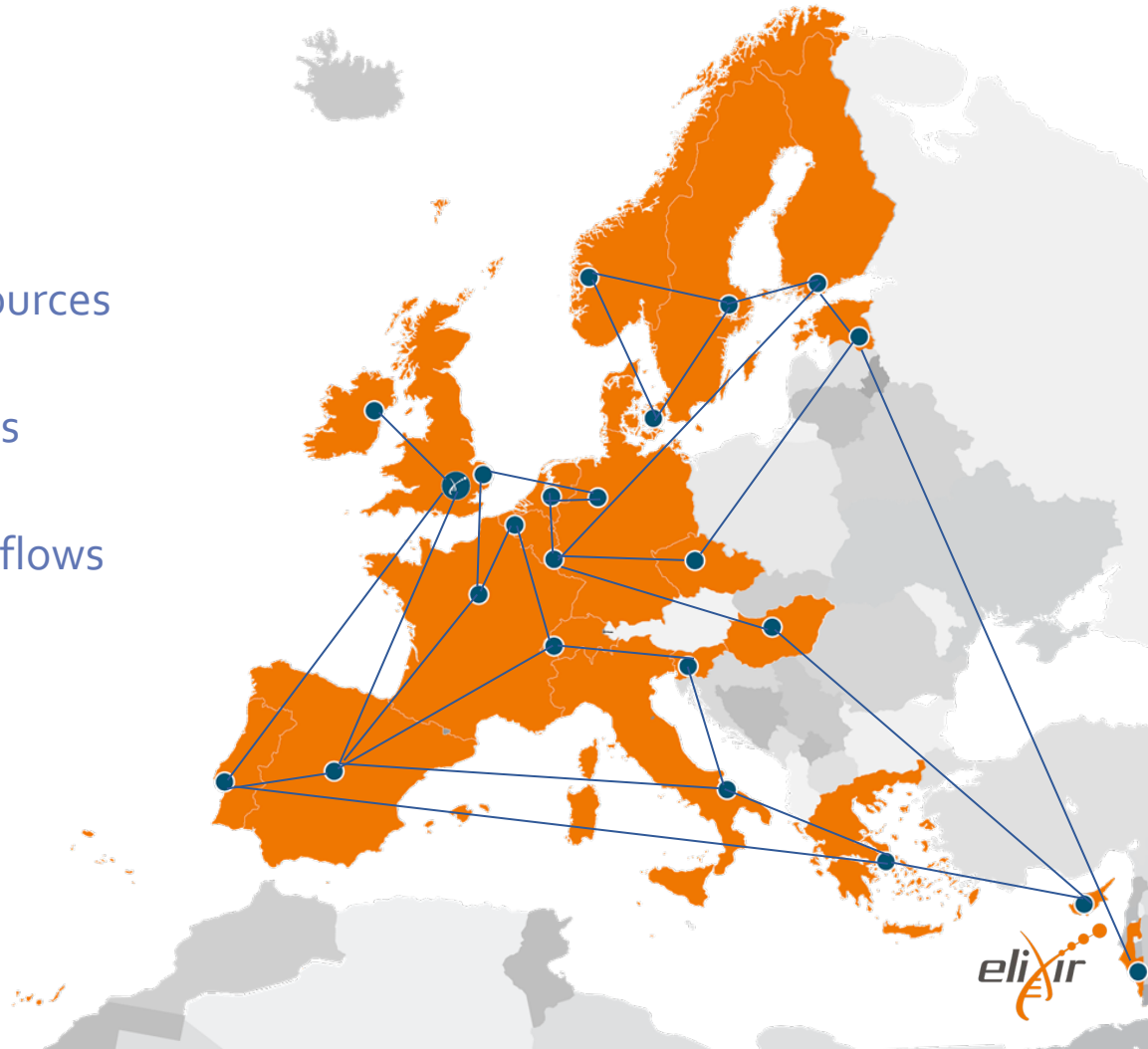
# Coordinating services

Databases and Data Resources

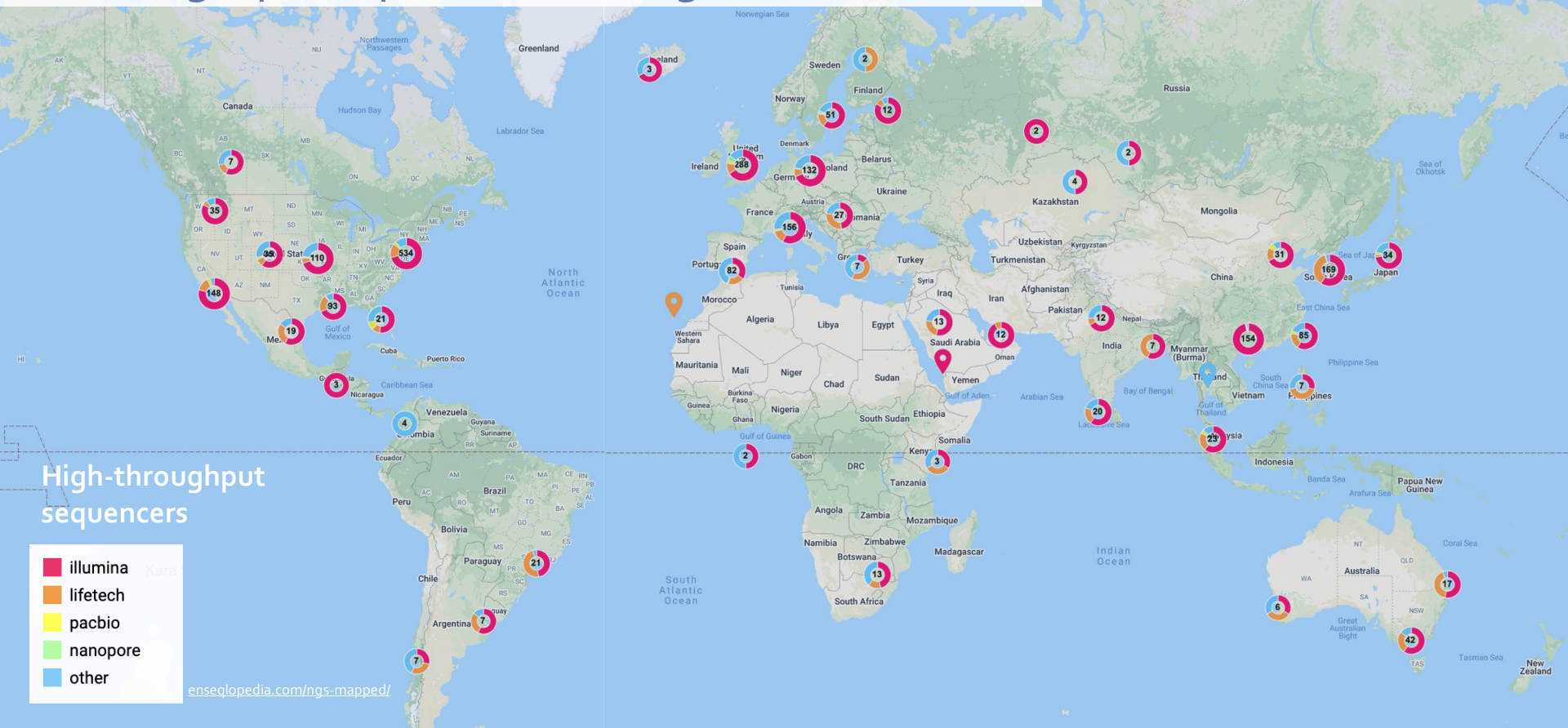Interoperability Resources

Tools, software and workflows

Compute Capabilities

Training Opportunities
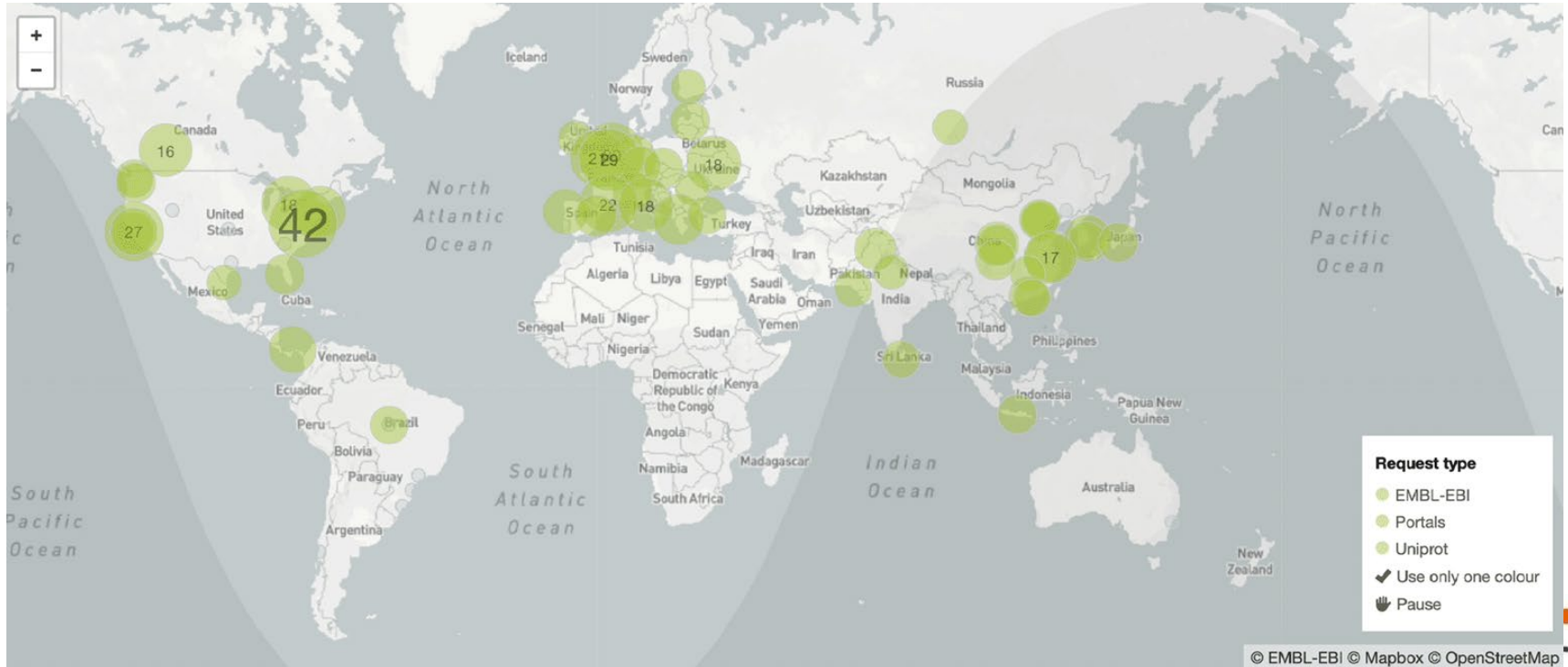
Data management

Geographic spread — data generation

High-throughput sequencers

- illumina
- lifetech
- pacbio
- nanopore
- other

enseqlopedia.com/ngs-mapped/

# Geographic spread — usage of data
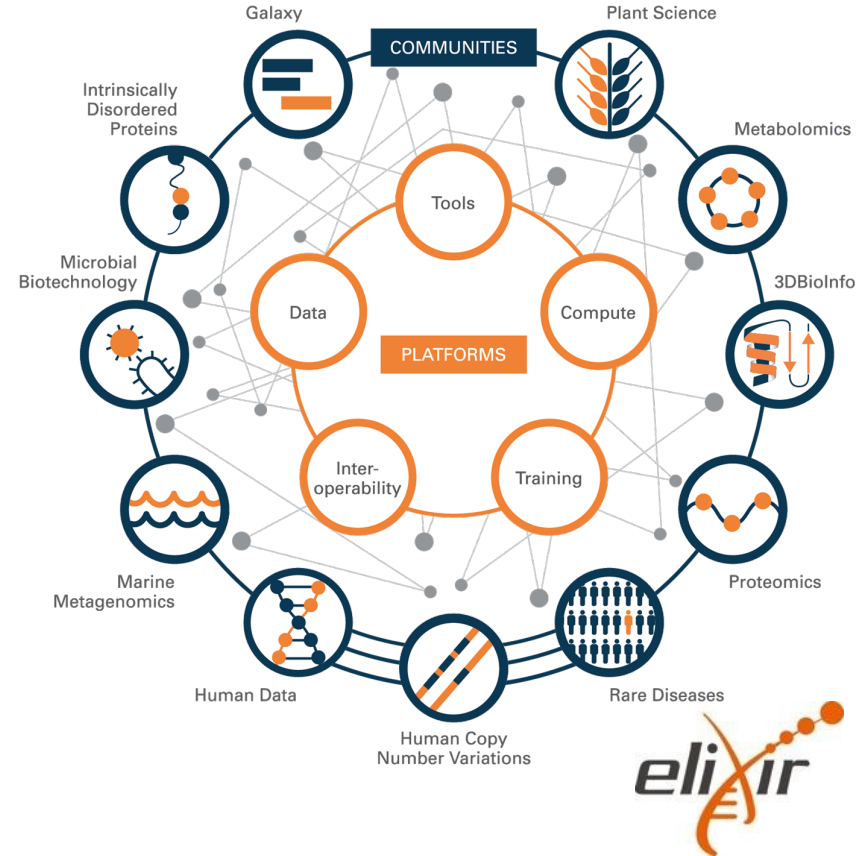## EMBL-EBI live data map

# ELIXIR and the Compute Platform

**ELIXIR Compute Platform**

- Identity and access management

- Data Integration for Compute

- ELIXIR Hybrid Cloud Ecosystem

- Deploying Reproducible Containers and Workflows across Cloud Environments

**ELIXIR Tools Platform**

- Packaging, Containerisation and Deployment

# Global Alliance for Genomics & Health

**Enabling responsible genomic data sharing for the benefit of human health**

The Global Alliance for Genomics and Health (GA4GH) is a policy-framing and technical standards-setting organization, seeking to enable responsible genomic data sharing within a human rights framework.

**Global Alliance for Genomics & Health**
Collaborate. Innovate. Accelerate.

650+ member institutions in 50+ countries

*"Collaborate on standards, compete on implementations!"*

- **Holistic approach** bringing together all stakeholders: patients, healthcare professionals, industry, academia, lawmakers...

- Developing modular, functional and interoperable **community standards & guidelines** driven by 24 Driver Projects

- Organized in **Work Streams** (e.g., Cloud, Data Security, Regulatory & Ethics)

- ELIXIR is a Strategic Partner of GA4GH

# How we work

**Represent** ELIXIR stakeholders in GA4GH & **promote** GA4GH standards within ELIXIR

**Prototype** real-world use cases with ELIXIR stakeholders, **develop** PoCs & **deploy** at ELIXIR nodes

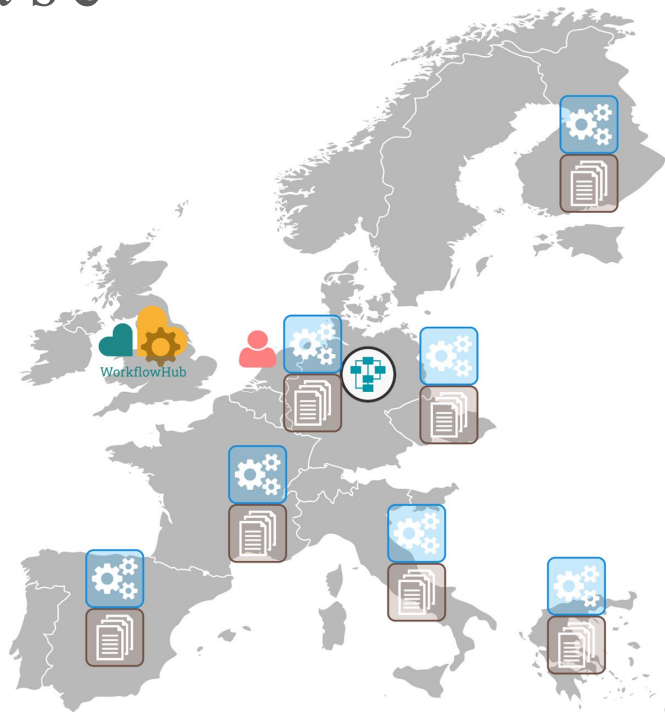**Consult** on integrating GA4GH standards into existing solutions and provide **technical support**

**Interoperability testing** with third party GA4GH-powered solutions

# An example use case



1. Process sensitive data across multiple ELIXIR nodes (**workflow** )

2. Collect processed data ( )

3. Centrally analyze at host institution (**workflow** )

# How to federate clouds?

Adoption of
- GA4GH standards
- Common governance & security guidelines

**Core services**
to further harmonize behavior & facilitate integration

**Reference implementations & auxiliary services**
to bridge gaps & facilitate uptake

# GA4GH Cloud and Access Management APIs

**Passport (AAI)**

*Grant access to data & compute*

**TRS: Tool Registry Service API**

*Access workflows and container images*

**DRS: Data Repository Service API**

*Access to data sets*

**WES: Workflow Execution Service API**

*Interpret workflows & schedule task execution*

**TES: Task Execution Service API**

*Execute tasks*

# Technological Implementations: Galaxy Pulsar



https://pulsar-network.readthedocs.io

# Galaxy Community becomes an EOSC Infrastructure...

- deferred data / remote data
- distributed storage
- scheduling taking storage into account

**HORIZON-INFRA-2021-EOSC-01 –
EuroScienceGateway just kicked off** ☺

# Workflow registry
# https://workflowhub.eu



## Workflows
- Any kind of system
- May remain in home repository
- Linked with data, documentation
- Can have DOIs

## WorkflowHub
- Links with other system, import/export (GA4GH TRS, RO-Crate)
- Metadata standards
- Github integration

## Mixed depth of support for Workflow Management Systems
- Lift metadata from different systems
- Coupled to execution platforms

# Coupled to execution environments

# ELIXIR::GA4GH enabled Services & Solutions

| Galaxy PROJECT | WorkflowHub | ELIXIR::GA4GH Cloud |
|---|---|---|
| Web-based platform for reproducible computational analysis | Registry for describing, sharing and publishing scientific computational workflows | Federated, interoperable network of workflow engines and compute nodes based on GA4GH standards |
| ELIXIR Community | EOSC-Life resource | GA4GH Driver Project |
| APIs & (third-party) GUIs | API & GUI | APIs & third-party GUIs |

Maturity

# National Implementations: de.NBI Cloud Federation



[https://cloud.denbi.de](https://cloud.denbi.de)

- Fully **academic cloud** federation

- Established 2016

- Provides **storage and computing resources** for the life sciences community

- **Free of charge** for academic use

- Federation is **maintained by the eight German cloud centers** located in Berlin, Bielefeld, Freiburg, Gießen, Heidelberg and Tübingen (+ FZ Jülich in 2022/2023)

- de.NBI Cloud offers a solution to enable **integrative analyses**, the **efficient use of data** in research, and computational **capacities for bioinformatics training**.

# WESKIT

**Features**
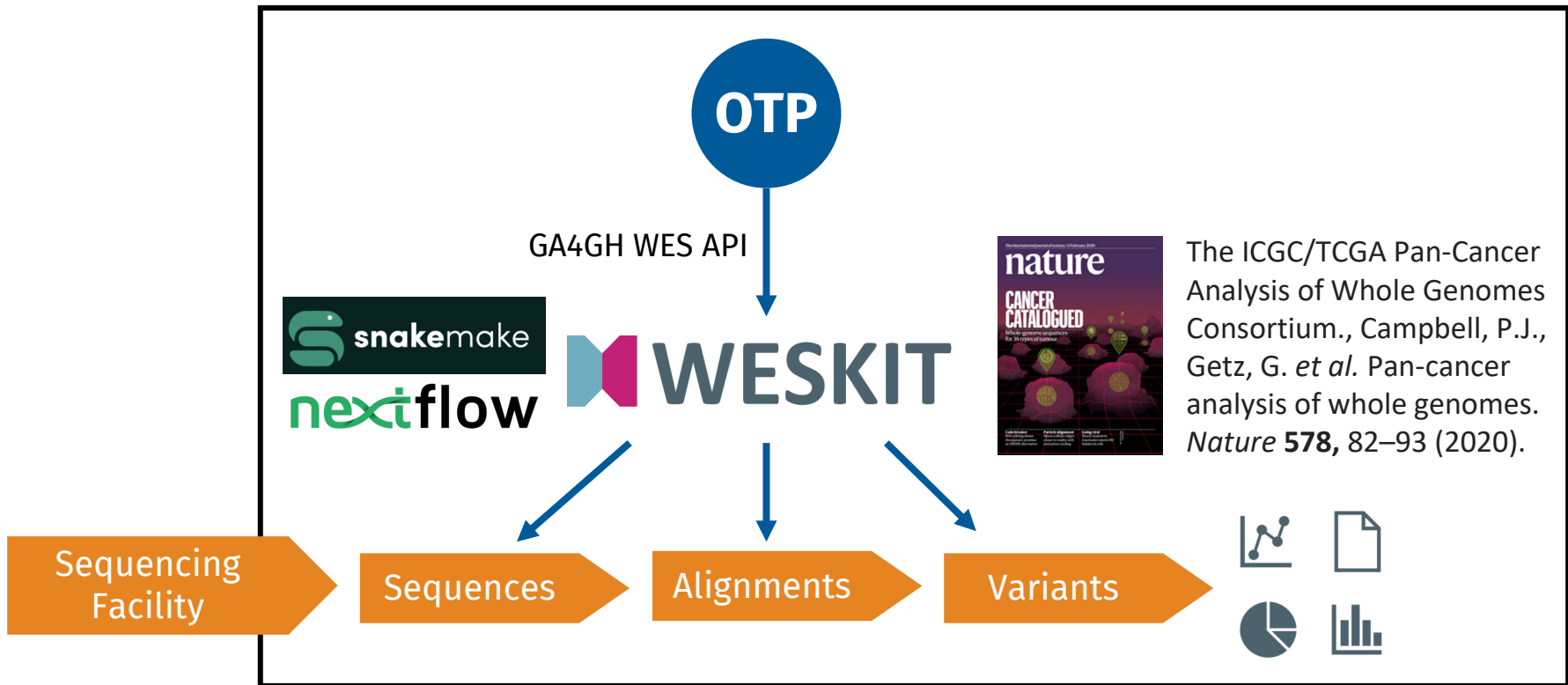- GA4GH WES implementation
- Execution of Snakemake and Nextflow workflows
- Focus: stability and high data throughput
- Developed and used at BIH, DKFZ & Sanger
- HPC and cloud deployment
- OIDC support
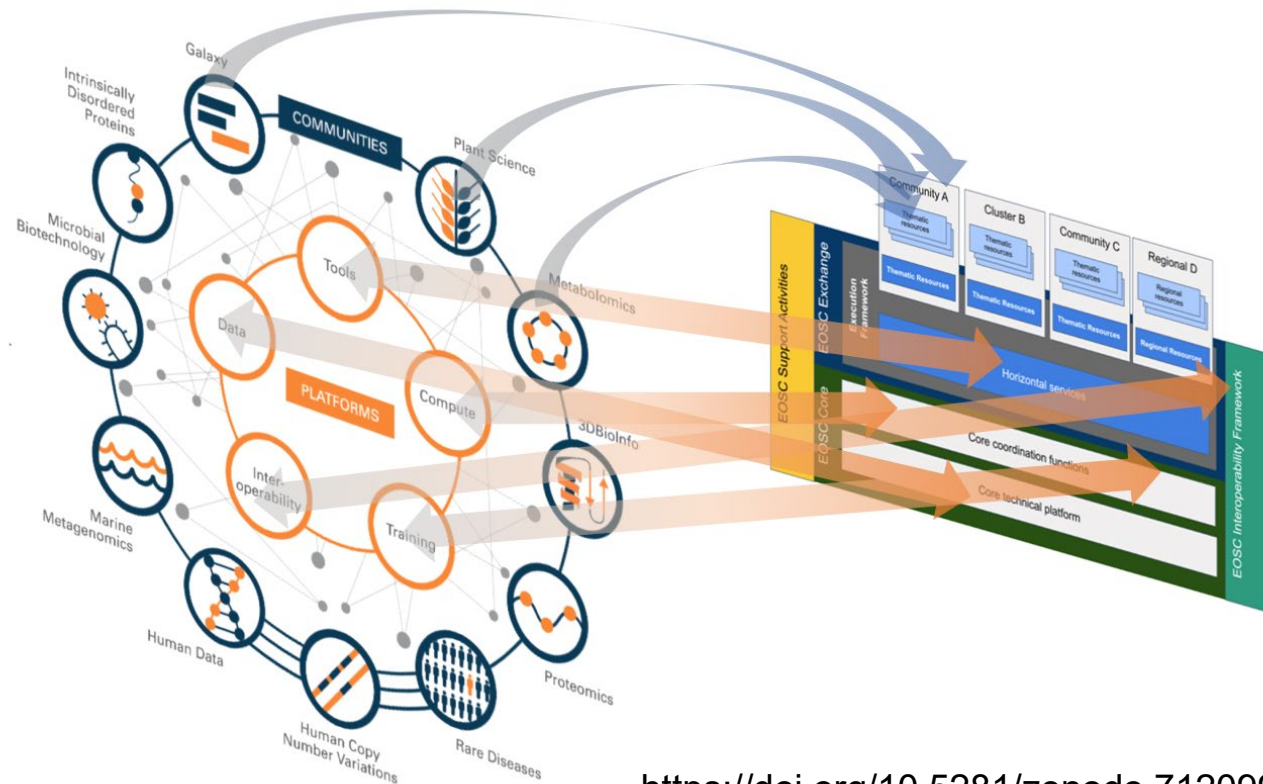- Developed in cooperation with ELIXIR compute platform

# ELIXIR's EOSC Strategy: ELIXIR is an existing Infrastructure - EOSC is a delivery partner

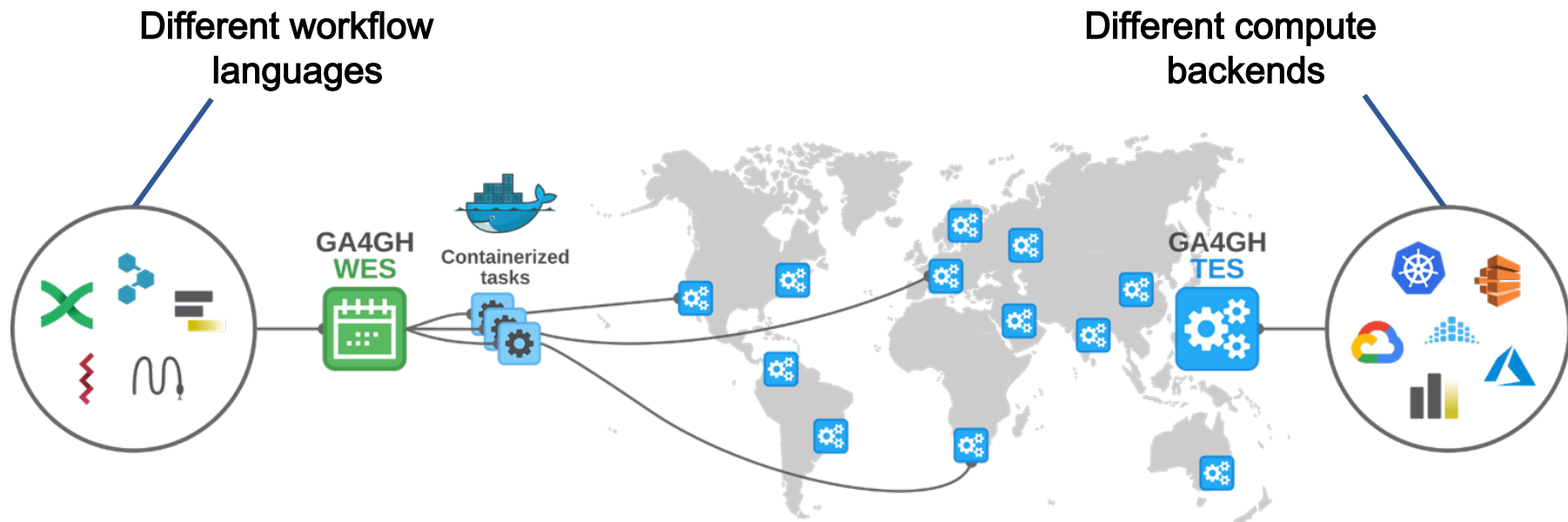ELIXIR's communities represent domain-specific expertise and resources.

ELIXIR's Commissioned Services represent an ecosystem which meets many of the same goals as EOSC.

We will co-develop EOSC, influence components where we have expertise and adopting technologies when they have demonstrated utility.



https://doi.org/10.5281/zenodo.7120996

# ELIXIR is enabling federated analytics through the use of GA4GH Cloud API specifications



Different workflow languages

Different compute backends

GA4GH WES

Containerized tasks

GA4GH TES

*"Send the compute to where the data is!"*

**Thanks**

*Contact: jonathan.tedds@elixir-europe.org*

*www.elixir-europe.org*