

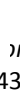
Collaborative Management Platform for  
detection and Analyses of (Re-) emerging  
and foodborne outbreaks in Europe

# A global platform for the sequence-based rapid identification of pathogens

Prof. Frank M. Aarestrup, coordinator (Technical University of Denmark)

Prof. Marion Koopmans, deputy coordinator (Erasmus Medical Center, the  
Netherlands)



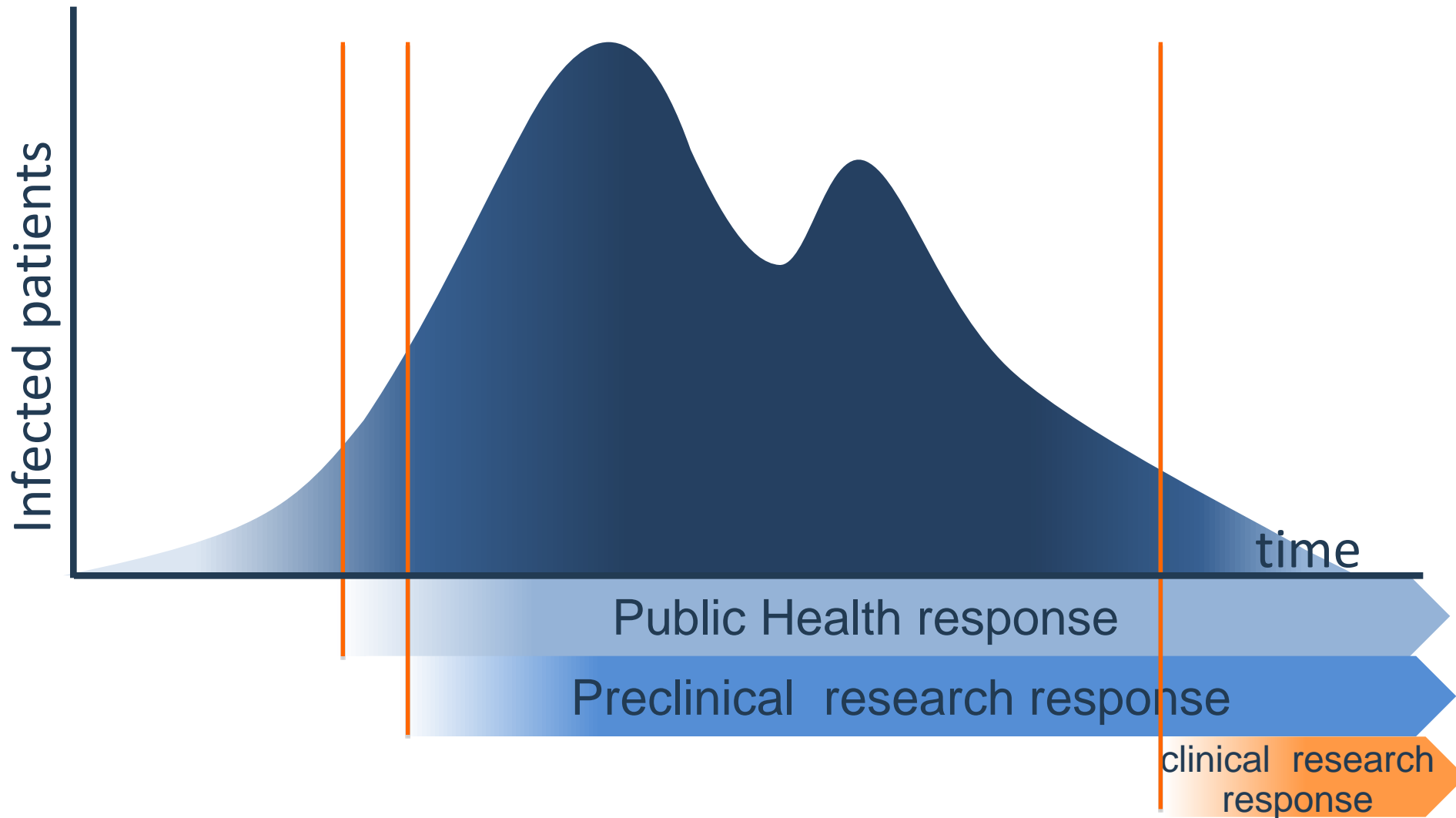
This project has received funding from the *European Union's Horizon 2020*  
*research and innovation programme* under grant agreement No 643476.   
*research and innovation programme* under grant agreement No 643



# Infectious disease situation 2015

- Dynamics of common infectious diseases are changing
  - Demographic change, population density, anti vaccine, AMR, etc.
- New diseases emerge frequently
  - Deforestation, population growth, health system inequalities, travel, trade, climate change
- Effects are difficult to predict due to complexity of problems
  - Rapid flexible response
- Public health and clinical response depend on global capacity for disease surveillance
  - Rapid sharing, comparison and analysis of data from multiple sources and using multiple methodologies

# Clinical research response to ID outbreaks usually fragmented and too late



# What the world needs

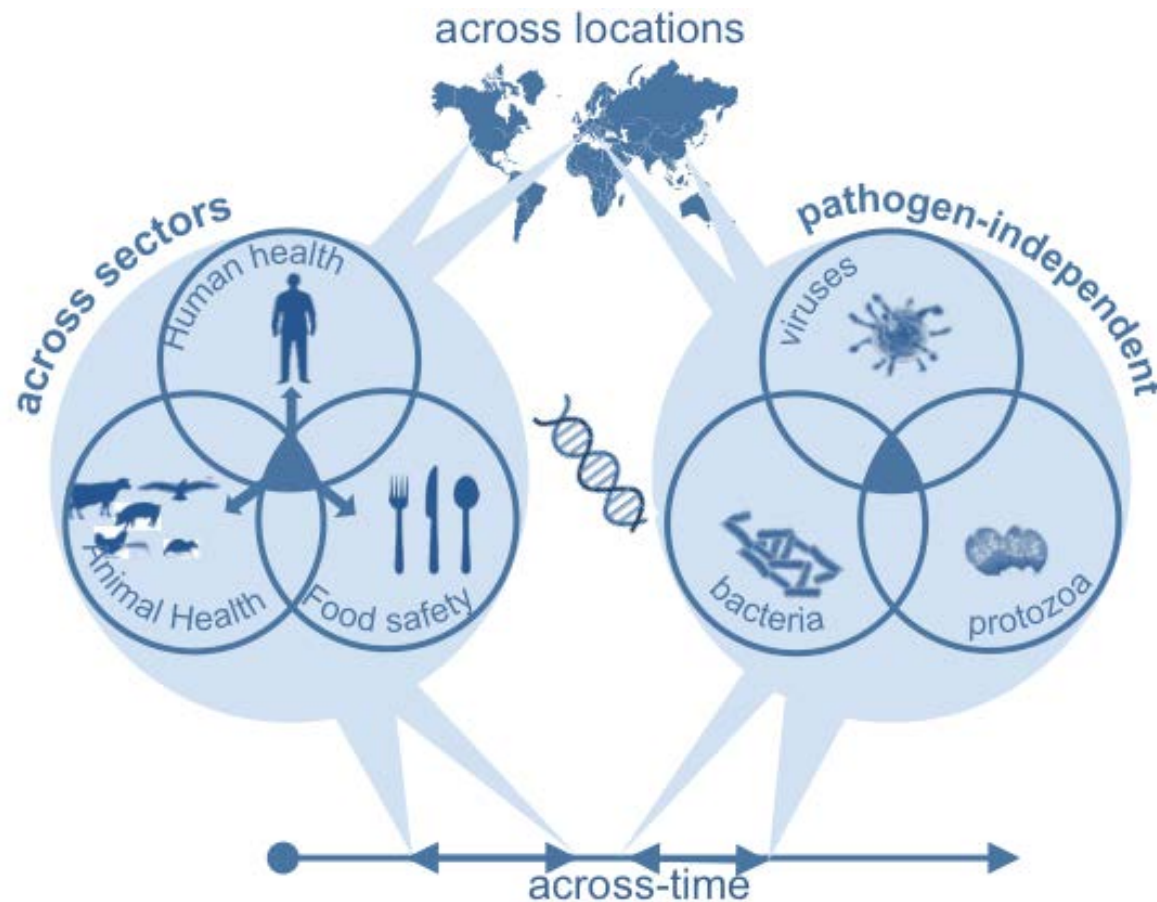
- Real-time data on occurrences of all infectious agents
- (Automatic) detection of related clusters in time and space
- Possibility to observe trends in clones and species as well as virulence and resistance
- Ability to rapidly compare between all types of data

**There can be no real-time surveillance without real-time data sharing**

# NGS advantages

- Laboratory diagnostics increasingly rely on (pathogen) genomic information
  - RNA / DNA are common across pathogens, therefore, methods to analyse pathogen genomes are potentially universal
  - Next generation sequencing capacity is developing fast, and costs are becoming competitive
- 
- Capturing NGS developments may provide a universal language that can be harnessed for early detection of outbreaks across disciplines and domains
  - If the technology keeps developing, less equipped labs may leapfrog

# Our vision: to build one system that serves all



# Epidemic preparedness research: European Union-supported efforts



- **Prediction**
  - Emergence, surveillance, modelling
- **Early recognition and containment**
  - Surveillance, clinical awareness, infection control
- **Data infrastructure**
  - Data repositories, sharing
- **clinical research**
  - pathogen & disease characterisation
  - prevention & treatment
- **funding**
  - rapid responses



2009-2016  
€ 36 M



2015-2020  
€ 21 M



2013-?  
€ >100 M



2014-2019  
€ 24 M



2015-2020  
€ 3 M

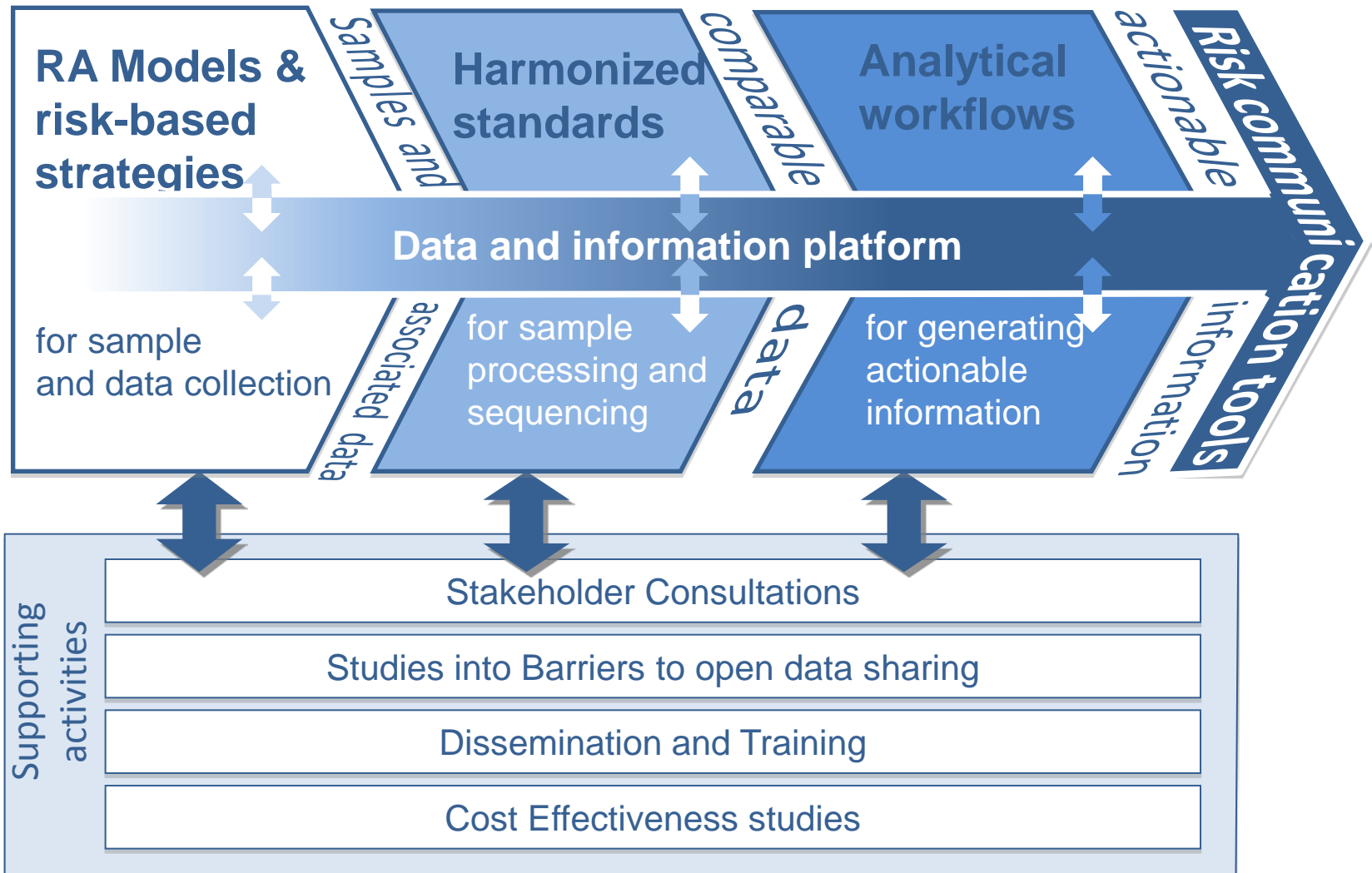


# COMPARE principles

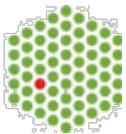
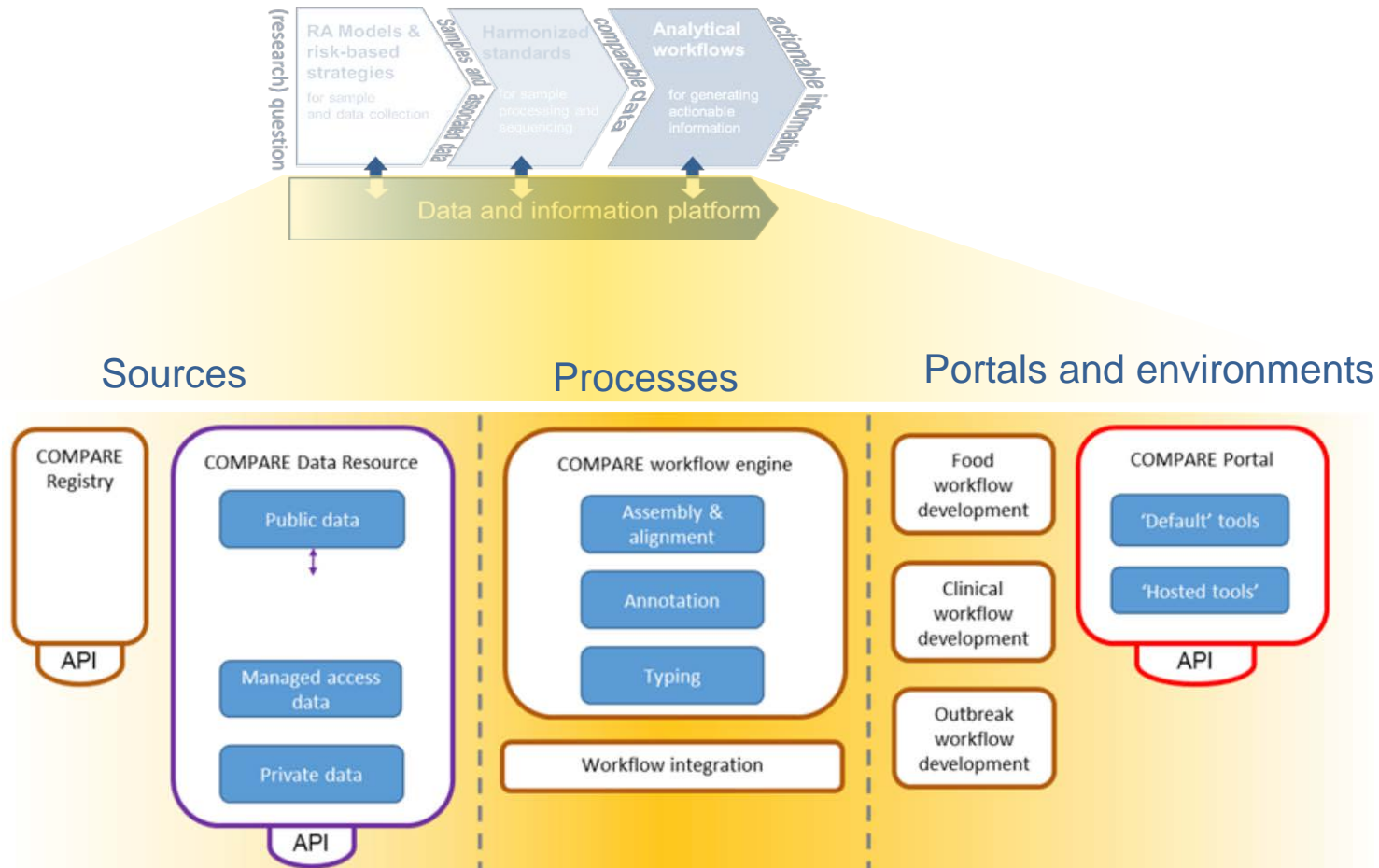
- COMPARE is sector-, domain- and pathogen-**independent**
- Analyzing **sequence-based pathogen data** in combination with **associated (clinical, epidemiological and other) data**
- Building on **established infrastructures (ENA, Elixir)**
- COMPARE is a **user-driven system**, designed with the information needs of its intended diverse group of future users and other stakeholders in mind
- COMPARE will make **optimal use of existing and future complementary systems, networks and databases**, ensuring compatibility where needed
- COMPARE is a **flexible, scalable and open-source based** information-sharing platform
- 1 December 2014 – 31 November 2019



# Project structure



# WP9 Information sharing platform

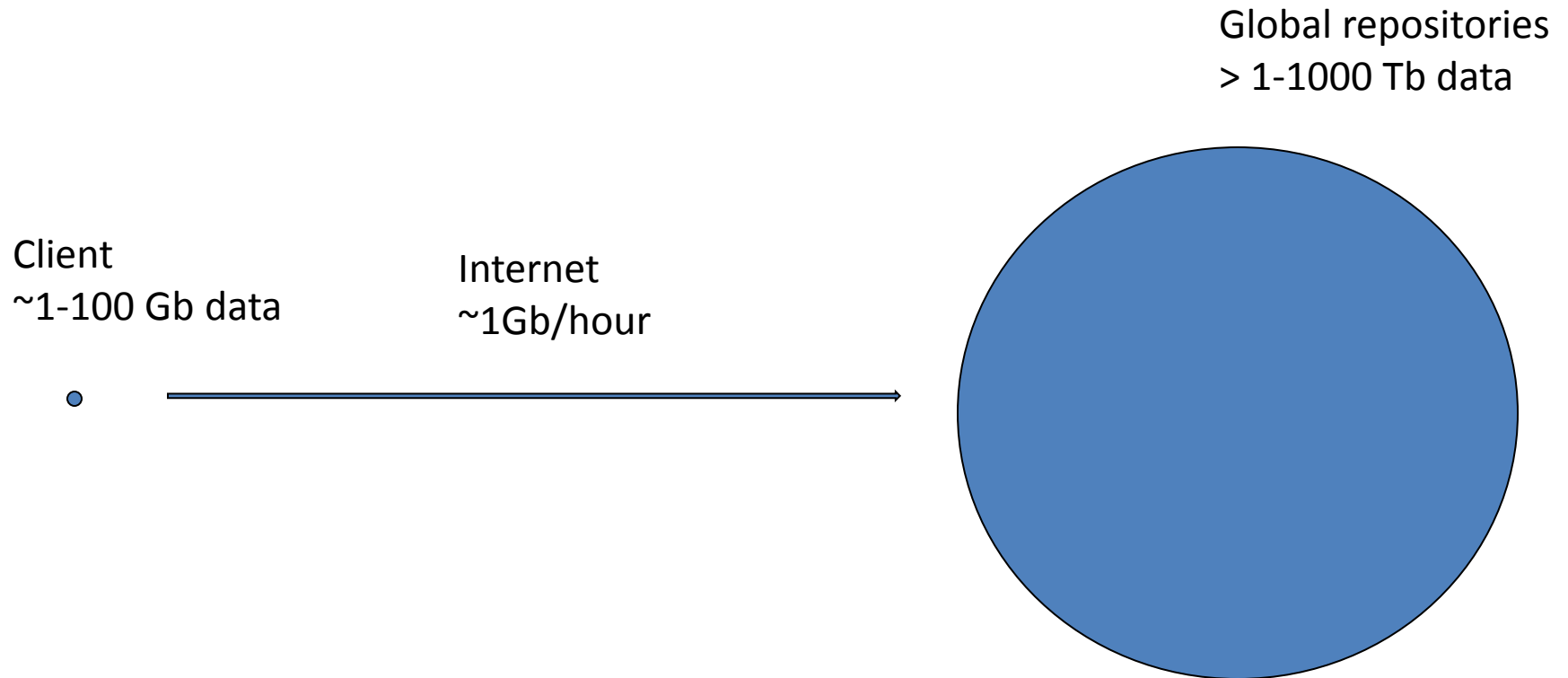


Building on the EU ESFRI Elixir, EMBL and DTU infrastructures

This project has received funding from the *European Union's Horizon 2020 research and innovation programme* under grant agreement No 643476.

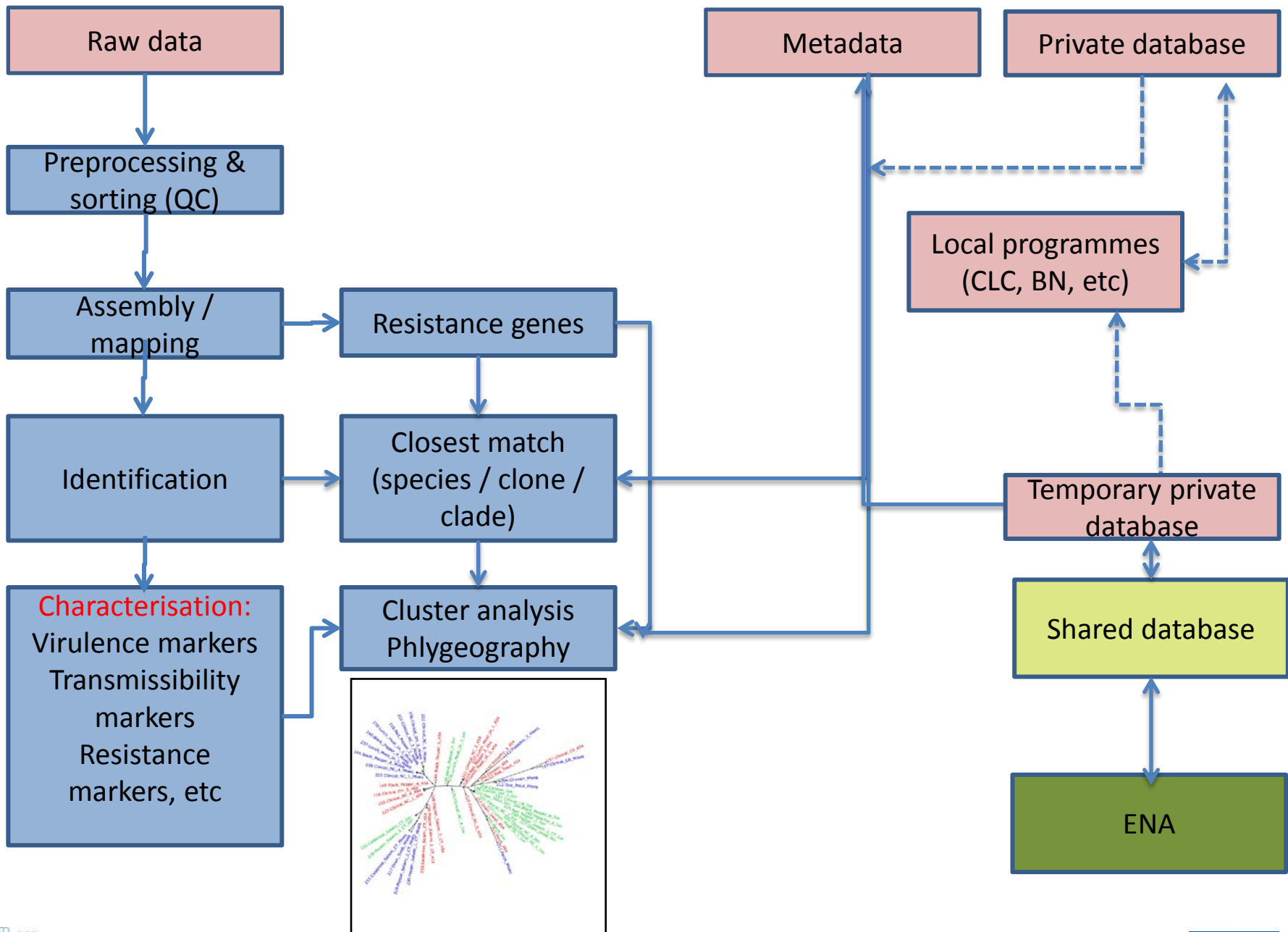


# Data comparison problem



# Workpackage 9: Data infrastructure design

- Sharing of highly structured and well-described data
  - COMPARE standards (e.g. checklists relating to isolates)
  - Data reporting tools that support the structuring and validation of data
- Support for a spectrum of data types
  - Raw NGS to assemblies
  - Derived data: typing information, AMR profiles, etc.
- Data availability
  - Early pre-publication private access, where required, for defined user groups according to explicit data-sharing agreements
  - Rapid flow of data to full public availability and global presentation
- Data discovery and retrieval systems, taking full advantage of data and metadata structures
  - Cloud-based autonomous analytical workflows (assembly, typing, phylogenetics, etc.)
  - Unifying data portal



# Current status - IT



The **COL**laborative **M**anagement **P**latform for detection and **A**nalyses of (**Re-**) emerging and foodborne outbreaks in Europe is a collaboration of 29 institutions with experience in outbreak detection and response in areas of human health, animal health and food safety.

## COMPARE Reference Genomes

This COMPARE Reference Genomes page offers a curated selection of published reference sequences covering viral, bacterial and protozoan genomes. These sequences can be searched and retrieved via the following URLs as tagged records in the European Nucleotide Archive (ENA). The complete COMPARE Reference Genomes dataset can be retrieved via the following URL:

<http://www.ebi.ac.uk/ena/data/xref/search?source=COMPARE-RefGenome>

ENA sequence or sample accessions for a single sample/isolate in the dataset can be returned using the following URL:

[http://www.ebi.ac.uk/ena/data/xref/search?source=<source>&source\\_accession=<source\\_accession>](http://www.ebi.ac.uk/ena/data/xref/search?source=<source>&source_accession=<source_accession>)

where `source_accession` is the isolate/sample name as shown in the table below, for example:

[http://www.ebi.ac.uk/ena/data/xref/search?source=COMPARE-RefGenome&source\\_accession=Beijing/55028/2007/CHN](http://www.ebi.ac.uk/ena/data/xref/search?source=COMPARE-RefGenome&source_accession=Beijing/55028/2007/CHN)

The ENA record, shown in the 'Target primary accession' column of the result from the above URL, can be retrieved with the following URL:

[http://www.ebi.ac.uk/ena/data/view/<Target\\_primary\\_accession>](http://www.ebi.ac.uk/ena/data/view/<Target_primary_accession>)

where `Target_primary_accession` has been inserted from the response to the previous URL (e.g.

<http://www.ebi.ac.uk/ena/data/view/GQ856465>).

More extensive functions are described for REST services relating to the COMPARE Reference Genomes. Users should note that records in the dataset are served from ENA and are denoted as belonging to the dataset through ENA cross-reference annotations.

The below table serves as a quick overview of the COMPARE Reference Genomes set with direct hyperlinks to the INSDC records. Text in brackets next to a taxon name represents serovar or genotype information.

Sample/Isolate ID	Aggregated taxonomic name or Taxon name	Genome	INSDC record(s)
Norwalk/1968/US	NoV/GI.P1/GI.1	complete	M87661
Southampton/1991/UK	NoV/GI.P2/GI.2	complete	L07418
NLV/VA98115/1998	NoV/GI.P3/GI.3	partial	AB287450
Chiba407/198/JP	NoV/GI.P4/GI.4	complete	AB042808
Musgrove/1989/UK	NoV/GI.P5/GI.5	partial	AJ277614
BSS/1997/DE	NoV/GI.P6/GI.6	complete	AF093797

# Reference genomes

- Curated selection of published reference sequences covering viruses, bacteria and protozoa
- Summary page and entry point at <http://www.ebi.ac.uk/ena/about/compare-reference-genomes>
- Searchable through webservice
- Launched 18.6.15, extended 25.8.15

## In progress

Protocols for sampling and handling

Ongoing benchmarking and ring trials



# COMPARE data hubs

- On-request service to share pre-publication data
  - Set up: provider, consumer accounts, scope, etc.
  - Providers report data through existing COMPARE interfaces
  - Consumers retrieve metadata (web or spreadsheet manifest) and data (Globus FTP)
- Several existing hubs, including
  - **dcc\_sibelius**, for the Influenza H5N8 pilot - pre-publication read data and metadata
  - **dcc\_compare**, for replication of overall COMPARE public content, e.g. for external clouds or other infrastructure; currently all bacteria, viruses and some parasites, in time refinements expected

Your Endpoint

globus Manage Data Groups Support

Transfer Files Activity Manage Endpoints Dashboard

Get Globus Core Turn your computer on

Transfer Files

Endpoint  Go

Path  Go

Please select an endpoint above.

Label This Transfer

© 2010-2016 Computation Institute, University of Chicago, Argonne National Laboratory

globus

Transfer F

Transfer Files

Endpoint  Go

Path  Go

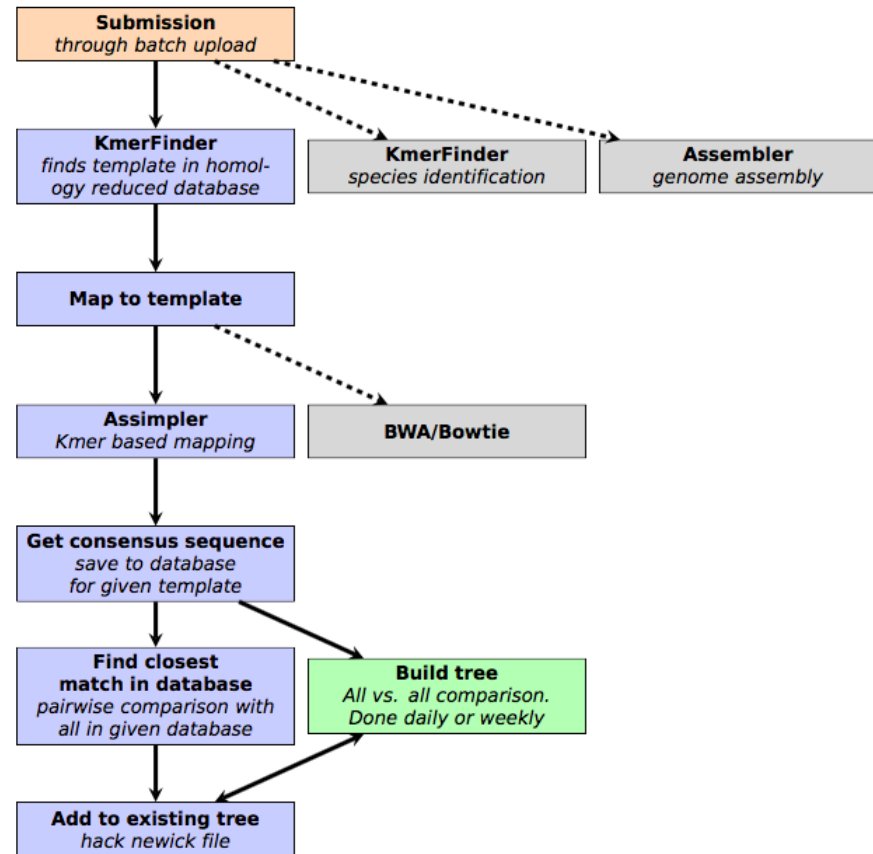
select all | none | up one folder | refresh list

- Folder dcc\_berlioz
- Folder dcc\_compare
- Folder dcc\_marc
- Folder dcc\_sibelius

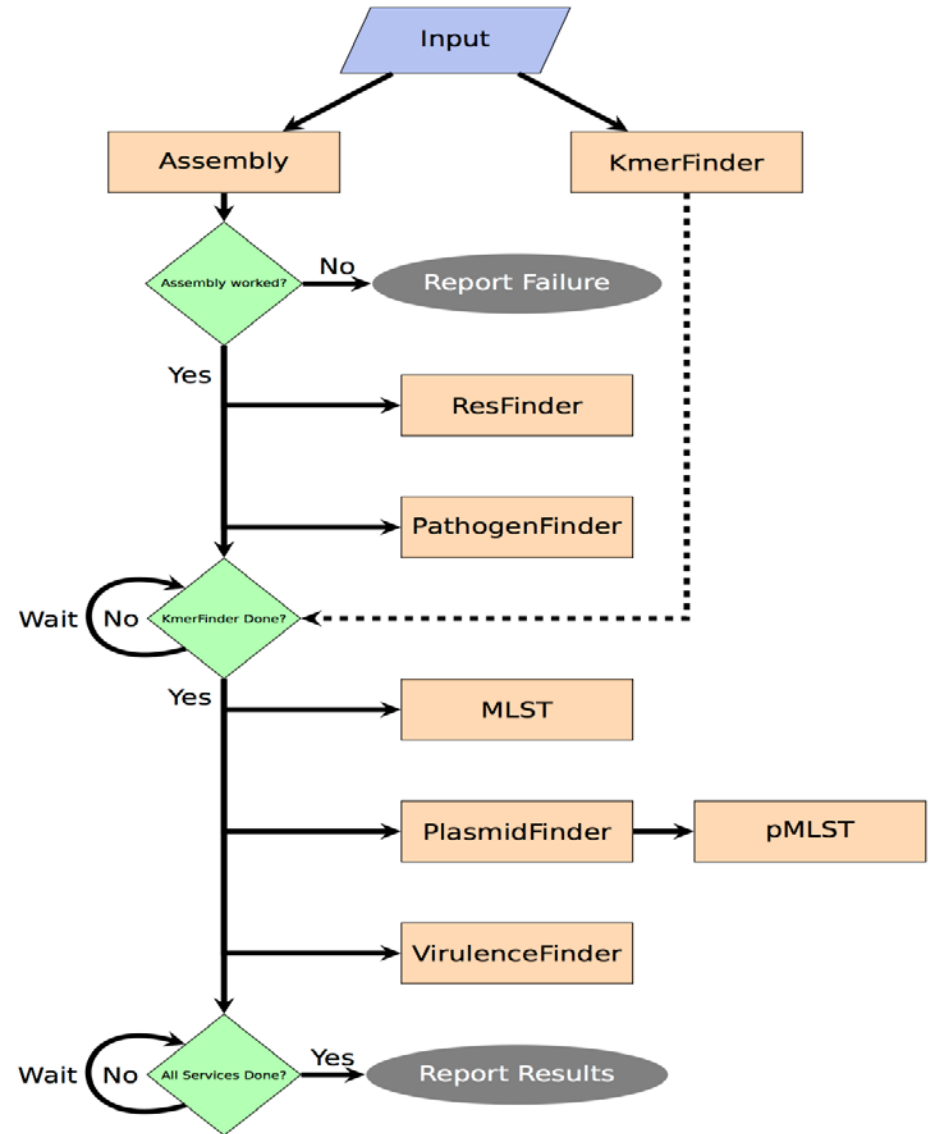


# Cloud compute

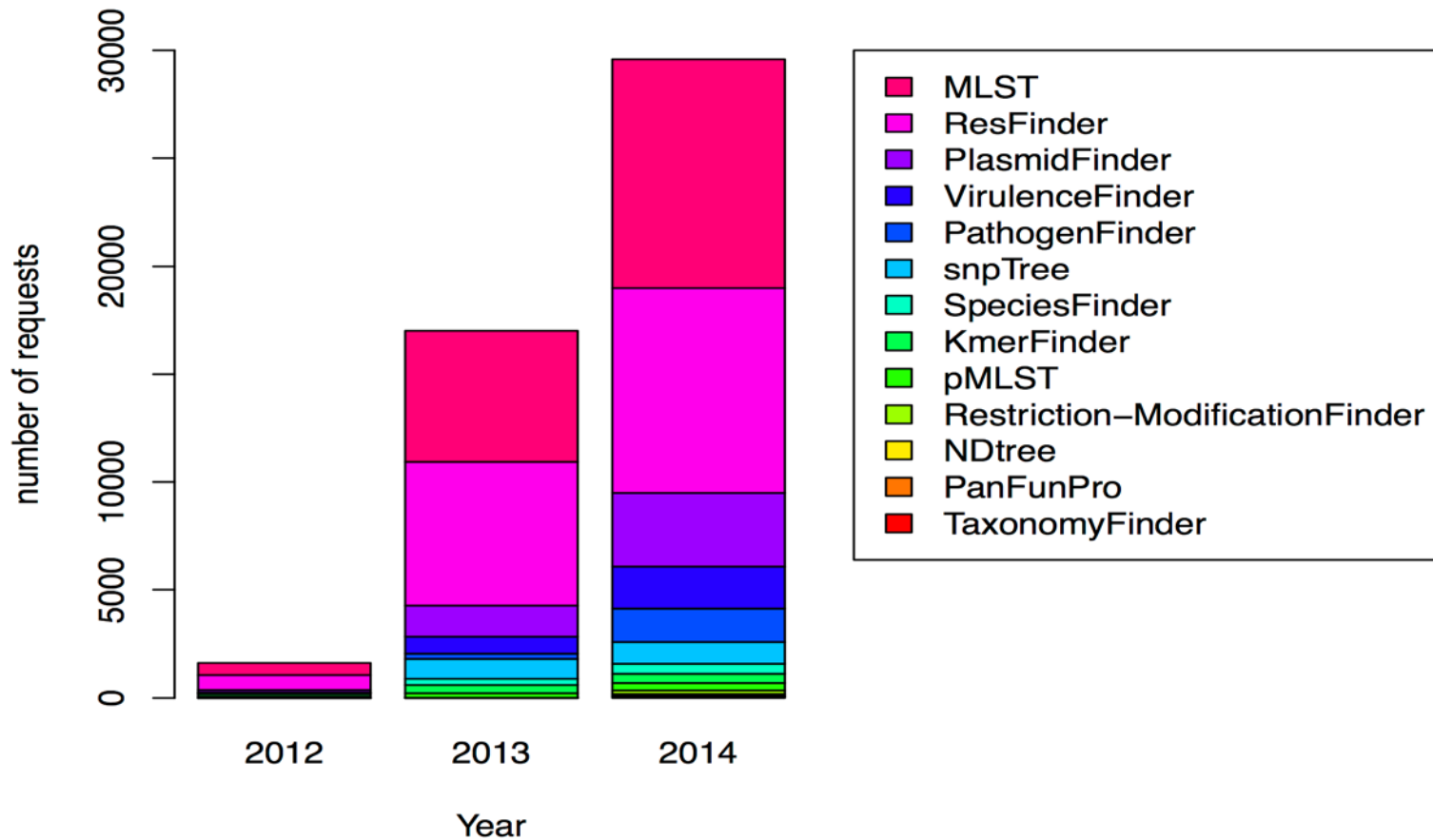
- Embassy environment
  - Available to COMPARE
    - OpenStack framework (in place)
    - BioLinux standard machine (in place)
  - In development
    - CGE Docker
    - Discussions with CLIMB, NECTAR, IFB-Cloud (Galaxy)
- Development of analysis workflows
  - Evergreen
  - Assembly (Velvet + Prokka)
  - Metagenomics
  - **CGE**
  - iPython environment for engagement of workflow environments



# Bacterial Analytic Pipeline



# User Statistics



Until now: >175,000 submissions

From + 8,000 IP-adresses

# Center for Genomic Epidemiology Welcome fmaa

- Home
- Services
- Instructions
- Example
- Article abstract

## Bacterial Analysis Pipeline - Batch Upload

The CGE Bacterial Analysis Pipeline executes a workflow of services with predefined parameters and stores the submitted data and result in the database at the user's disposal.

View the [version history](#) of this server.

**STEP 1:** [Download Metadata Template](#) [Template](#)

**STEP 2:** Fill out template

**STEP 3:** [Upload Metadata File](#)

**STEP 4:** [Select Files](#)

**STEP 5:** [Submit](#)

### Supported browsers

Browser	Internet Explorer	Firefox	Chrome	Opera	Safari
≥ V 8.0					
V 10.0	✓	✓			
edge	✓	✓			
≥ V 36.0	✓	✓	✓	✓	✓
≥ V 41.0	✓	✓	✓	✓	✓

### Progress Overview

Name	Size	Progress	Status

[Remove all](#)

- Support
- Scientific problems
- Technical problems

Copyright DTU 2011 / All rights reserved  
Center for Genomic Epidemiology, DTU, Kemitorvet, Building 204, 2800 Kgs. Lyngby, Denmark  
Funded by: The Danish Council for Strategic Research  
Last modified August 25, 2015 13:32:18 GMT



This project has received funding from the *European Union's Horizon 2020 research and innovation programme* under grant agreement No 643476.



EN\_A\_upload\_template.xlsx - Microsoft Excel

The insert size must be given as an integer (i.e. no decimals)

Field Name	Mandatory	Description
sample_name	No	The name of the isolate the user uses to identify the sample
group_name	Not Yet Used	The name of the group. Adding a group name will incorporate the sample in the analysis pipeline execution of that group
file_names	Yes	Name of all files associated to this sample. Multiple filenames should be separated by a space
sequencing_platform	Yes	Choose between: LS454, Illumina, Ion Torrent, ABI SOLiD and unknown
sequencing_type	Yes	Choose between: single, paired, mate-paired, and unknown
pre_assembled	Yes	Has the uploaded sample data been assembled? yes / no
sample_type	Yes	Choose between: isolate or metagenomic
organism	Yes	Write 'unknown' if the organism name is unknown
strain	No	Strain type ID
subtype	No	Subspecific genetic lineage, i.e MLST, serovar and biotype
country	Yes	
region	No	
city	(No)	
zip_code	(No)	Please provide as much information as possible. Low resolution locations reduce the usability. The minimal recommended option is to provide either city, zip_code or longitude and latitude coordinates.
longitude	(No)	
latitude	(No)	
location_note	No	Additional relevant details about the location
isolation_source	Yes	The host from which the sample/isolate has been taken. This should be a proper scientific name or one of the phrases found in the host cheat of this document.
source_note	No	Additional relevant details about the isolation source. i.e. blood, laboratory experiment or urine
pathogenic	Yes	Is the organism decreed pathogenic? yes / no / unknown
pathogenicity_note	No	Additional relevant details about the organism's pathogenicity
collection_date	Yes	The date of the sample collection. Use one of the following format: YYYY-MM-DD or YYYY-MM or YYYY
collected_by	No	Name of the institute or person who took the sample
usage_restrictions	Yes	Choose either 'private' or 'public'. Note that private data will be deleted after some time to free disk space on the server.
release_date	Not Yet Used	Write the date from when the data and results should be public available. Format: YYYY-MM-DD
email_address	Not Yet Used	Email adress of the uploader
notes	No	Any additional information can be added
insert_size	Yes	The insert size must be given as an integer (i.e. no decimals)

Ready

14:17 03-09-2015



This project has received funding from the *European Union's Horizon 2020 research and innovation programme* under grant agreement No 643476.



CGE Server <https://cge.cbs.dtu.dk/services/cge/> Welcome fmaa

# Center for Genomic Epidemiology

Example Article abstract

eters and stores the submitted data and result in the database at the users

Open

Computer > Local Disk (C:) > Search Local Disk (C:)

Organize New folder

File list:

Name	Date modified	Type
Copy of OCS 162 - WGS_Trivl Rqst - REGISTRATIO...	03-11-2015 16:15	Microsof
Copy of Oplæg_2016_frank.xlsx	05-10-2015 10:57	Microsof
Copy of Metonym Budget DTU_v00.xlsx	04-09-2015 08:43	Microsof
Copy of Copy of ENA_upload_check.xlsx	27-08-2015 08:50	Microsof
ENA_upload_check.xlsx	26-08-2015 15:38	Microsof
Copy of ENA_upload_template.xlsx	24-08-2015 09:21	Microsof
Copy of BAP-test-VTEC_raw_trimmed.xlsx	20-08-2015 08:17	Microsof
Copy of BAP-test-VTEC_raw_8-10.xlsx	19-08-2015 13:47	Microsof
Copy of BAP-test-VTEC_raw_8-12.xlsx	19-08-2015 13:00	Microsof
Copy of BAP-test-VTEC_raw_3_7.xlsx	19-08-2015 09:46	Microsof
Copy of BAP-test-VTEC_raw_3-12.xlsx	19-08-2015 08:56	Microsof
Copy of BAP-test-VTEC_raw_13-24.xlsx	19-08-2015 08:56	Microsof
Copy of BAP-test-VTEC_raw_38-47.xlsx	19-08-2015 08:55	Microsof
Copy of BAP-test-VTEC_raw_25-37.xlsx	19-08-2015 08:55	Microsof
Copy of BAP-test-VTEC_raw_25-47.xlsx	19-08-2015 08:21	Microsof
Copy of BAP-test-VTEC_raw_3-24.xlsx	19-08-2015 08:21	Microsof
Copy of BAP-test-VTEC_raw_onesandtwo.xlsx	14-08-2015 16:06	Microsof
Copy of BAP-test-VTEC_raw.xlsx	14-08-2015 14:34	Microsof

Select a file to preview.

File name: Microsoft Excel Worksheet

Open Cancel

Supported browsers

Browser	Windows	Mac OS X	Linux	Ubuntu
≥ V 8.0				✓
V 10.0	✓	✓		
edge	✓	✓		
≥ V 36.0	✓	✓	✓	✓
≥ V 41.0	✓	✓	✓	✓

Name	Size	Progress	Status
Remove all			

Support Scientific problems Technical problems

Copyright DTU 2011 / All rights reserved  
 Center for Genomic Epidemiology, DTU, Kemitorvet, Building 204, 2800 Kgs. Lyngby, Denmark  
 Funded by: The Danish Council for Strategic Research  
 Last modified August 25, 2015 13:32:18 GMT

11:41 16-11-2015



This project has received funding from the *European Union's Horizon 2020 research and innovation programme* under grant agreement No 643476.



Sample Overview x [https://cge.cbs.dtu.dk/tools\\_new/client/platform/sample/](https://cge.cbs.dtu.dk/tools_new/client/platform/sample/)

# Center for Genomic Epidemiology

Welcome fmaa

[Home](#)    [Services](#)    [Batch Upload](#)    [MapViewer](#)

## Sample Overview

[Download all data in an Excel spreadsheet?](#)  
[Download the resistance data in a more detailed Excel spreadsheet format?](#)

	Name	Date	Country	City	Origin	Action
-	VTEC47	2012-01-01	Denmark		human	Analyse Edit Remove
	<b>Service</b>	<b>Date</b>	<b>Status</b>	<b>Action</b>		
	<a href="#">KmerFinder-2.1</a>	2015-08-27	Success	Remove		
	<a href="#">Assembler-1.0</a>	2015-08-27	Success	Remove		
	<a href="#">ResFinder-2.1</a>	2015-08-27	Success	Remove		
	<a href="#">VirulenceFinder-1.2</a>	2015-08-27	Success	Remove		
	<a href="#">ContigAnalyzer-1.0</a>	2015-08-27	Success	Remove		
	<a href="#">PlasmidFinder-1.2</a>	2015-08-27	Success	Remove		
	<a href="#">MLST-1.6</a>	2015-08-27	Success	Remove		
	<a href="#">pMLST-1.4</a>	2015-08-27	Success	Remove		
+	VTEC46	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC45	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC44	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC43	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC42	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC41	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC40	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC39	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC38	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC37	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC36	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC35	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC34	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC33	2012-01-01	Denmark		human	Analyse Edit Remove
+	VTEC29	2012-01-01	Denmark		human	Analyse Edit Remove



This project has received funding from the *European Union's Horizon 2020 research and innovation programme* under grant agreement No 643476.







# Current developments - IT

- Ability to create own shared site
- Combining map, phylogeny and analysis (microreact)
- Evergreen trees
- Bring your own tools (already ability for own DB)

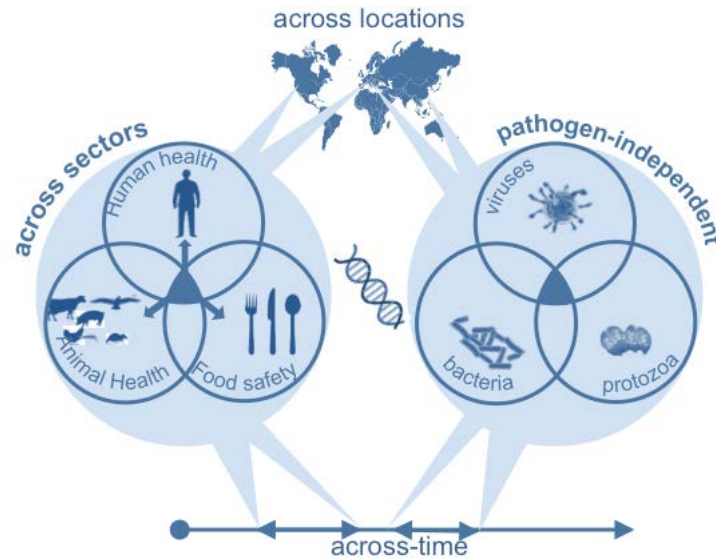
# Conclusions

- WGS/NGS is rapidly entering diagnostic and public health arena, with near real-time data generation
- Sequence platforms rapidly developing, cheaper, simpler
- Bottleneck at level of bioinformatics, particularly for intergroup comparison, national, international
- COMPARE aims to develop infrastructure and ICT to meet the coming demand
- In the coming years, we will be seeking partners for pilot projects

# Pilot projects in 2016

- Cross cutting: *Escherichia coli*, ESBL, norovirus, metagenomics
- Clinical WP: AMR
- Public health WP: 6 food-borne pathogens
- Emerging infections WP: Influenza and MERS CoV
- Ad hoc

# Our vision: one system serves all



## Guiding principles:

- Cross sector, cross domain, open source (not commercial)
- Interaction with the rest of the world (all inclusive)
- Data for action (actionable outputs)
- Central repository (ENA, DDJ, NCBI) (bring the tools to the data)

**There can be no real-time disease detection & surveillance without real-time data sharing**