



Data Sharing

Speaker: Kjell Petersen & Korbinian Bösl (UiB)

Moderator: Nazeefa Fatima (UiO)



Federated EGA Norway - a controlled access archival service for FAIR sharing of Norwegian human genome and phenotype data

Kjell Petersen

Technical Coordinator ELIXIR Norway



Learning Objectives

In this session, we will learn:

- the main principles of a controlled access archive
- how it is possible to also share sensitive data in life sciences
 - ...even if data sharing is tightly regulated
- how ELIXIR Federated Human Data sharing services supports FAIR data sharing



Session Take-Away

After completing this session, you will be:

- familiar with some key challenges and solutions to sharing of sensitive data
- able to understand how to get started



Learning Activity

Optional 10 min pair session
if you find the time (tomorrow)

Group activity; study the data submission and retrieval guidelines on <https://ega.elixir.no>. Discuss the steps researcher A need to go through to submit data to FEGA Norway and researcher B to retrieve the same dataset.

Data management

Data life cycle



Your role

Your domain

Your tasks

Data management

Data life cycle

Your role

Your domain

Your tasks

Compliance monitoring

Data analysis

Data management plan

Data organisation

Data protection

Data publication

Data quality

Data storage

Data transfer

Documentation and metadata

Existing data

Identifiers

Licensing

Machine actionability

Sensitive data

Your tasks

Sensitive data

- Is your data sensitive?
- How can you de-identify your data?
- Related pages
- More information
- Relevant tools and resources

Is your data sensitive?

Description

In general, the term "sensitive data" is used for any data that could do harm (for example to people, organisations, countries, or ecosystems) if it would be openly available. This can for example be personal commercial information, but also information such as breeding grounds of endangered species. Any data that is considered sensitive must be protected against unauthorized access. What is considered sensitive information is usually defined by national laws and may differ between countries. You should be cautious when you are dealing with sensitive, or potentially sensitive, information.

Considerations

- If you deal with any information about individuals from the EU, you are bound by the [General Data Protection Regulation \(GDPR\)](#). In GDPR, such data is called "personal data".
- In the context of GDPR "special category data" is a subclass of "personal data" that is potentially harmful, and GDPR prescribes very strict rules for dealing with this data. Article 9 of GDPR defines special categories as data consisting of racial or ethnic origin, political opinions, religious or philosophical beliefs, or trade union membership, genetic data, biometric data, data concerning health or data concerning a natural person's sex life or sexual orientation. Confusingly, these special categories are sometimes colloquially called "sensitive data". Note that this page is concerned with the broader concept of "sensitive data".

The way science work - and a scientific community

- Building on the work of others
 - “standing on the shoulders of giants” - or many, many, many researchers before me
- Life Sciences is a data driven/centric science field
- Sharing of data is essential in several phases
 - Among collaborators - in the active project phase
 - To the community and public - when publishing results, findings and interpretations
- Sensitive data such as human genomic data and phenotype information
 - Can not be shared open publicly for anyone to see
 - Controlled Access is required



Controlled Access Archival Service

- Dataset existence and description made public available, but not the data itself
- Access to data requires approval by authority controlling the dataset
 - Data Access Committee
- Download and reuse of data may require signed agreements



Overview of data management

FAIR data management from raw data provider to public repository

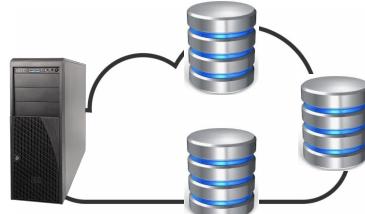
Data production



Data storage/sharing



Data compute



Data archiving



NTNU | HUNT Cloud

Federated European Genome-phenome Archive Norway

elixir NORWAY



Data for life



Sensitive data

GDPR

§

Personvernforordningen

Personopplysningsloven

Helseforskningsloven

NOR
MEN



Examples of sensitive data

Genomic sequences

5' TGACATGGGTACACATGACGGG 3'
||| ||| ||| ||| ||| ||| ||| |||
3' ACTGTACCCATGTGTACTGCC 5'



Disclaimer

I do not have:

- any formal legal education
- authorization or official role mandated from UiB

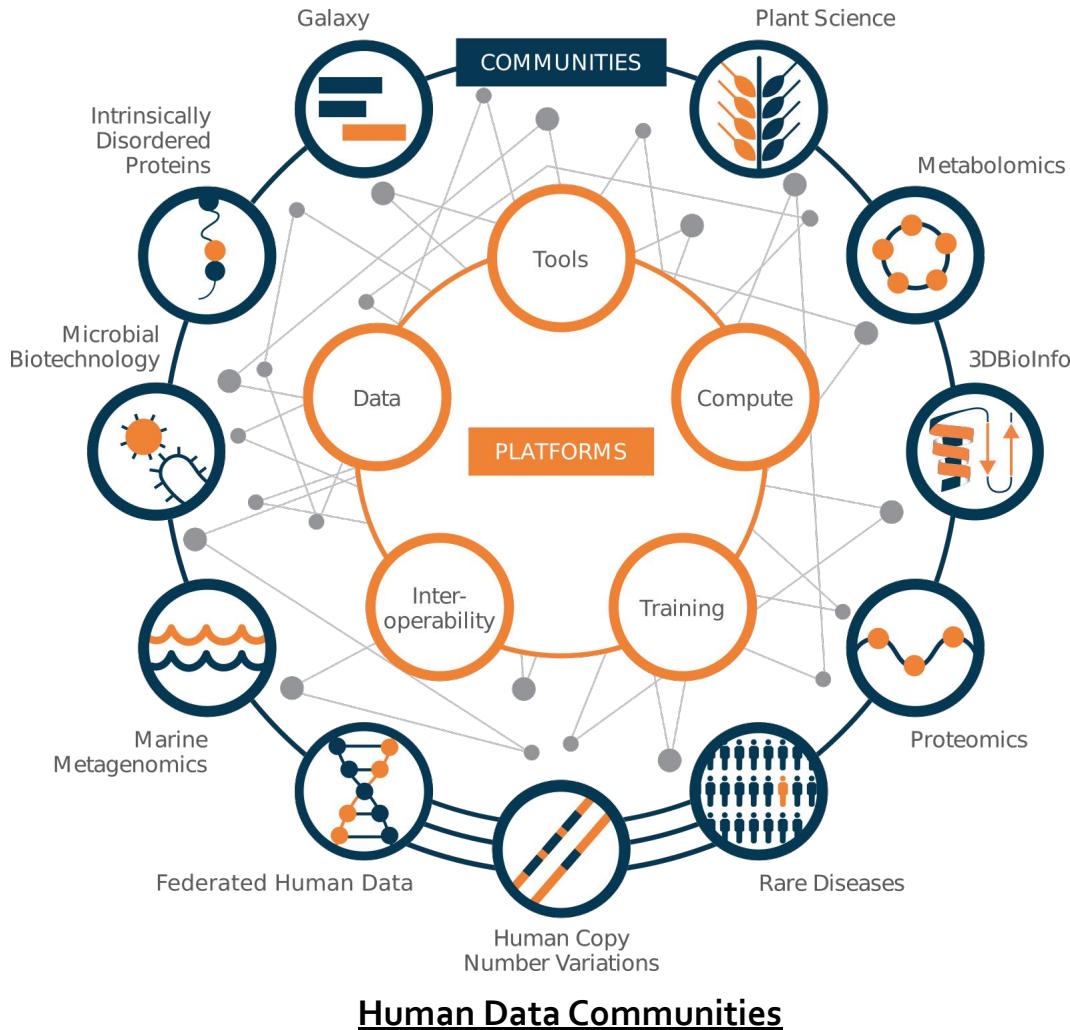


Content overview

- Brief background on ELIXIR Europe and international context for sharing sensitive data
 - Standards and international initiatives



ELIXIR Structure



Cancer Data Focus Group

- tumour heterogeneity;
- investigation of new enzymes;
- the fight against contamination;
- tumour evolution;
- cancer progression;
- harmonization of data from electronic health records.

**15 ELIXIR
Nodes
involved**

Health Data Focus Group

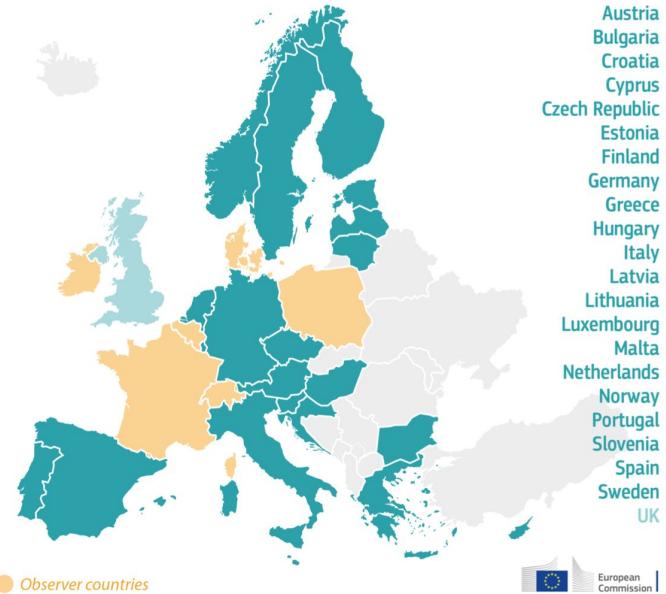
- omics
- clinical (structured and unstructured EHRs, and clinical imaging)
- lifestyle/patient questionnaire
- administrative
- real-world wearables/mobile devices data.



1+MG Declaration of cooperation - April 2018



EU countries agreed to cooperate in linking genomic data across borders



22 countries have now signed; 6 are observers

Aligning with international initiative



Working towards GA4GH standards, APIs and toolkits to be used throughout ELIXIR Nodes for human data discovery and access –
GA4GH into Europe



ELIXIR coordinating data federation

Minimum recommendations for EU-wide infrastructure to access and analyse genomic data

Standards

Application Programming Interfaces

Secure, federated clouds

Tools, services & workflows

Necessary minimal infrastructure component	In development ^a	Implemented at scale ^a
Genomics data and clinical information standards geared towards specific disease communities	Yes	No
Common application programming interfaces to enable remote data discovery and access	Yes	Yes
Computational resources, including secure, federated cloud computing environments that offer secure access across national boundaries to raw data and interoperable results	Yes	Yes
Regulatory frameworks that enable access to and the processing of genomic data across borders, including the management of transnational user access and compliance	Yes	No
A repository of tools and services, including workflows to analyse deposited data while enabling these analysis workflows to operate on data across national borders. This will contribute towards data reproducibility and provenance, which are of high importance in both research and clinical practices	Yes	Yes
A training and capacity-building programme to develop the skills and workforce required for genomics and big data in health care as well as shift the culture towards openness and integration of research data across national boundaries	Yes	Yes

^a'In development' and 'Implemented at scale' refer to locally defined status within ELIXIR and/or Biobanking and Biomolecular Resources Research Infrastructure.

[Saunders G et al \(2019\)](#)

nature REVIEWS GENETICS

Roadmap | Published: 27 August 2019

Leveraging European infrastructures to access 1 million human genomes by 2022

Gary Saunders, Michael Baudis, [...] Serena Scollon

Nature Reviews Genetics (2019) | Download Citation

Abstract

Human genomics is undergoing a step change from being a predominantly research-driven activity to one driven through health care as many countries in Europe now have nascent precision medicine programmes. To maximize the value of the genomic data generated, these data will need to be shared between institutions and across countries. In recognition of this challenge, 21 European countries recently signed a declaration to transactionally share data on at least 1 million human genomes by 2022. In this Roadmap, we identify the challenges of data sharing across borders and demonstrate that European research infrastructures are well-positioned to support the rapid implementation of widespread genomic data access.

Legal and regulatory frameworks

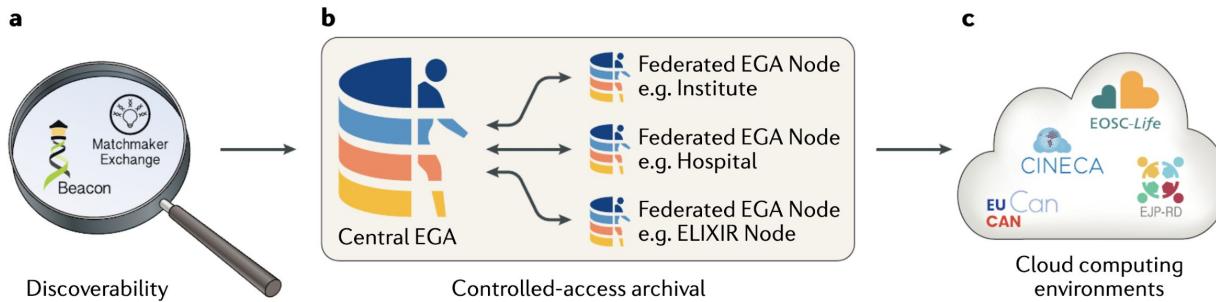
Training

Capacity building programme



Developing standards and driving use

Driving interoperability using GA4GH standards



EGA brief history

- The EGA is a resource for permanent secure archiving and sharing of all types of potentially identifiable human genetic and phenotypic data
- Implementing the FAIR principles

How the EGA is managed

The EGA was launched in 2008 by the EBI

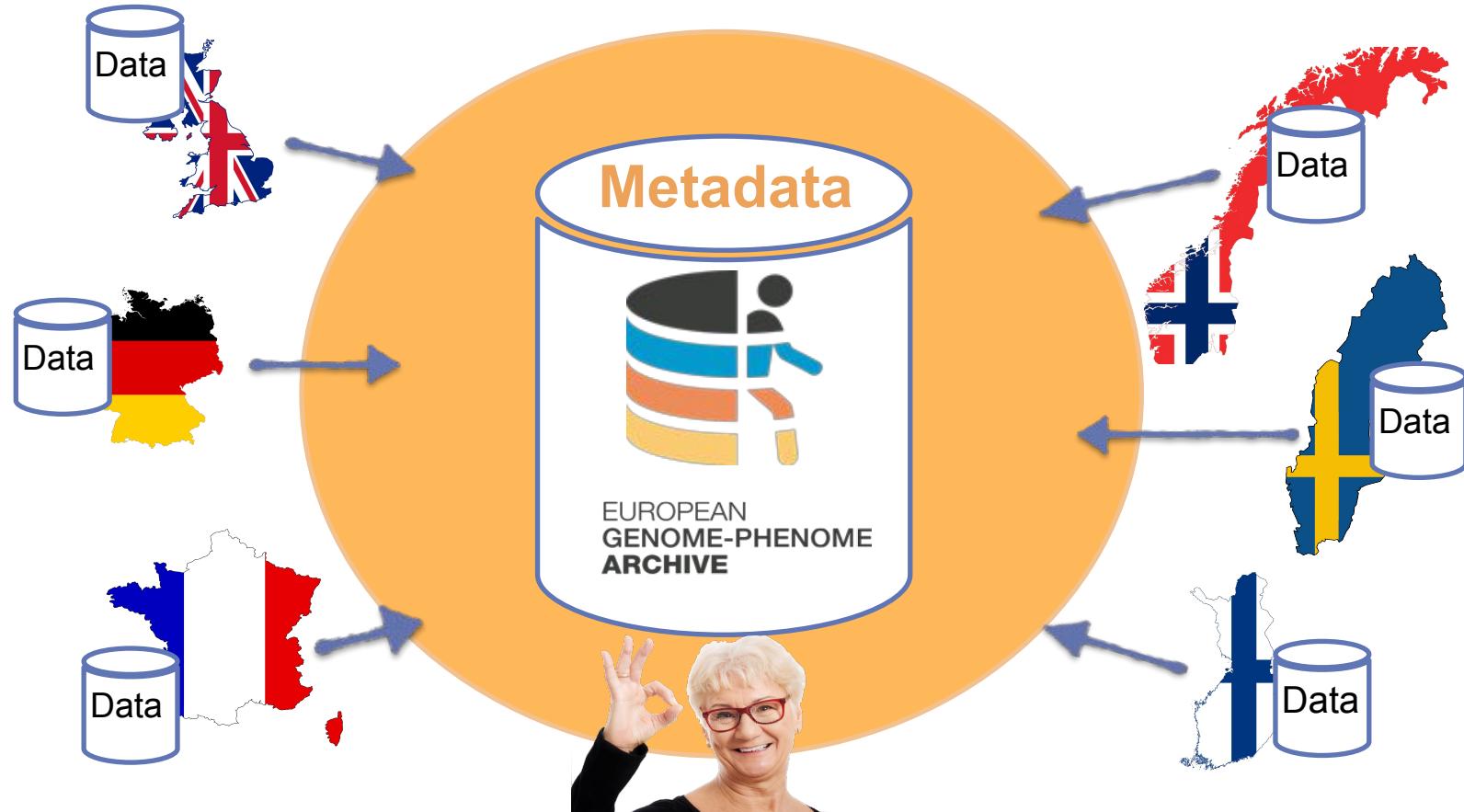


Why - and how is this FAIR?

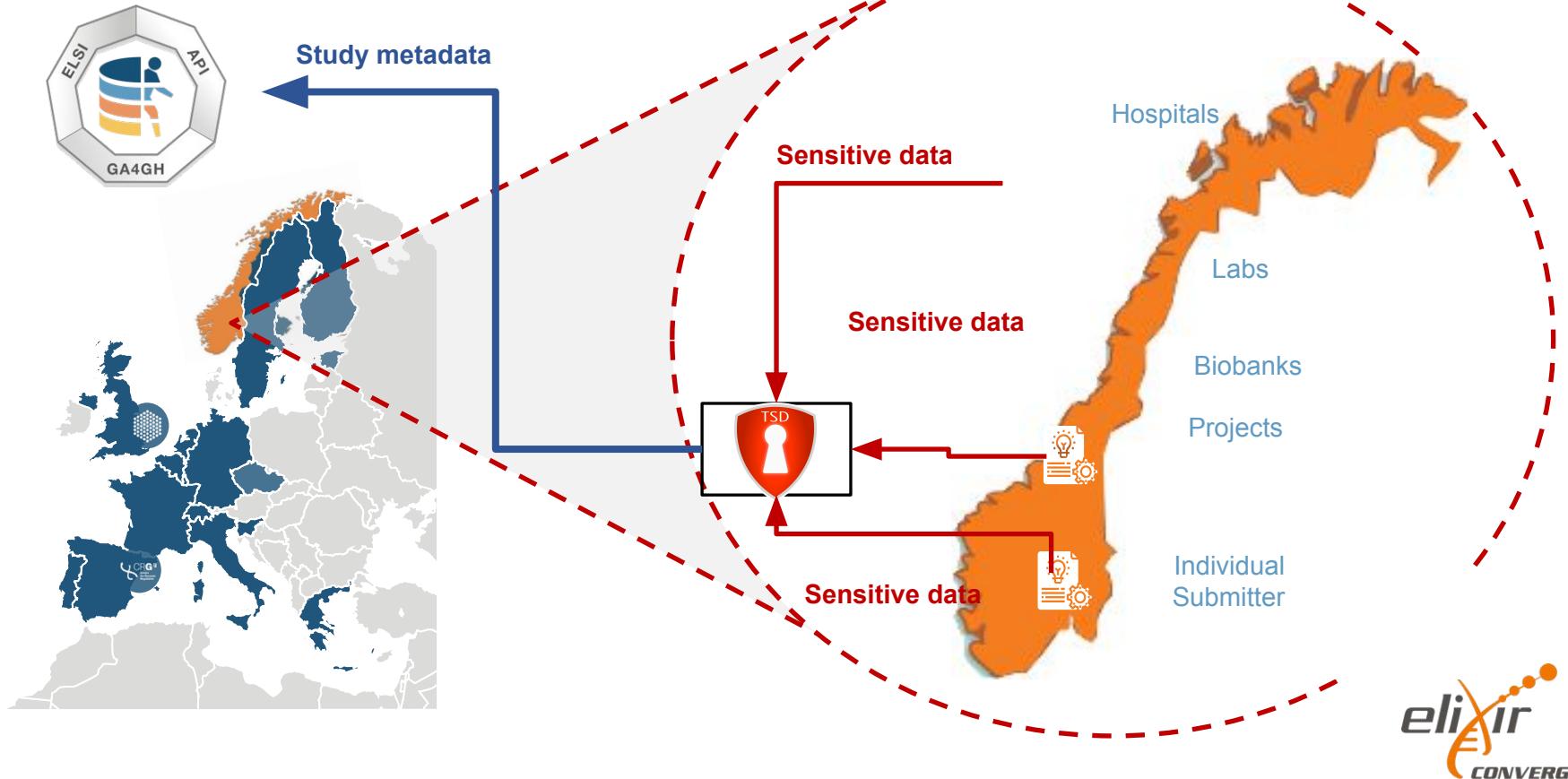
F : Findable

- Accession ID in publications to EGA / Federated EGA
- Your datasets existence are discoverable from the web portal

Federated EGA makes EGA scalable and GDPR-compliant



Building data federation capacity at European nodes



Some key GDPR concepts

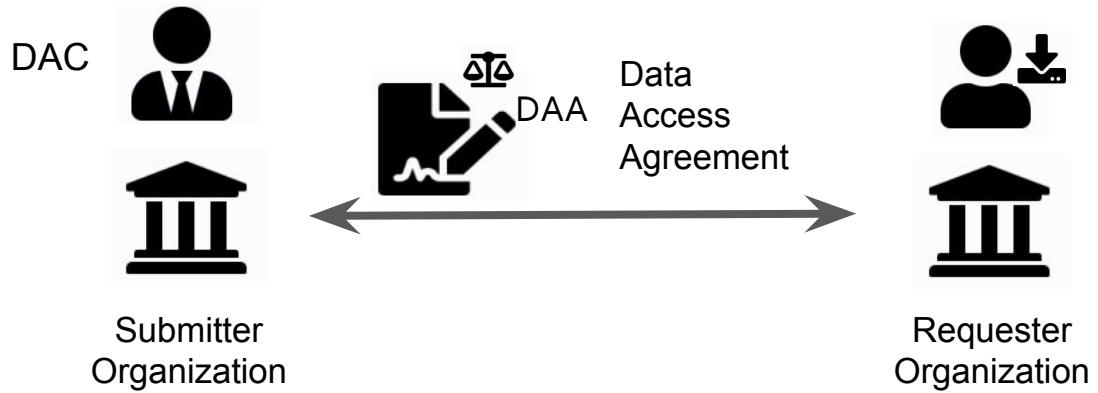


Data
Service

Data Controller

Data Processor

Data Access Agreement when sharing data





[Docs and Info](#) [Privacy Policy](#) [Terms of Service](#) [Contact Us](#) [About](#)

Federated EGA Norway node

 LS LOGIN

The Federated EGA (European Genome-phenome Archive) is a distributed solution for sharing and exchange of human -omics data across national borders. Federated instance collects metadata of -omics data collections stored in national or regional archives and make them available for search through the main EGA portal.



Searchable

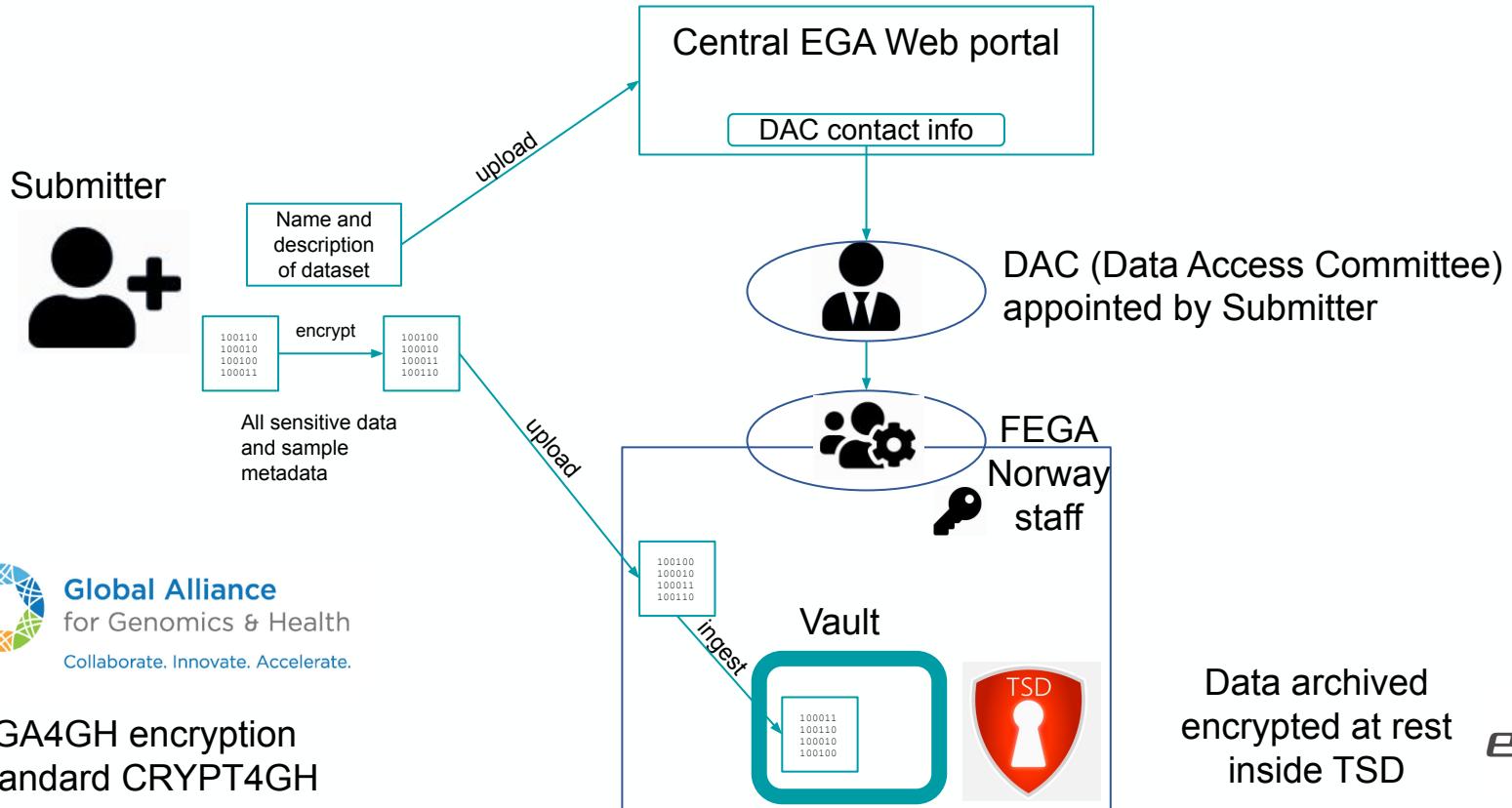


Secure



Shareable

Data flow and data control - part I





Search...

[Tips on how to search](#)[ABOUT](#) [SUBMISSION](#) [BROWSE](#) [ACCESS](#) [DOWNLOAD](#) [METADATA](#)[Helpdesk](#)[Log in](#)

The European Genome-phenome Archive (EGA) is a service for permanent archiving and sharing of all types of personally identifiable genetic and phenotypic data resulting from biomedical research projects.

DISEASE TYPE

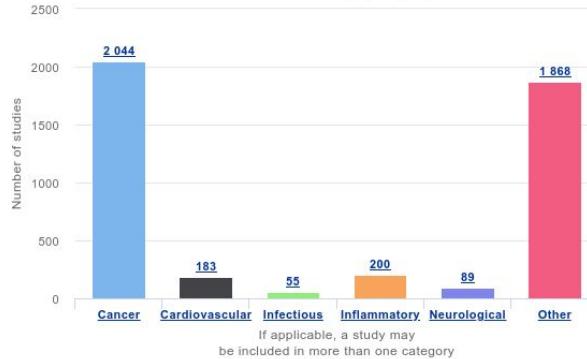
TECHNOLOGY

SAMPLE TYPES

What is in the EGA?

Studies in the EGA by disease

Click on a column to view category subgroups



Latest studies

The growing need for controlled data access models in clinical proteomics and metabolomics. *– 2021-10-01*

Telomere maintenance by telomerase activation or alternative lengthening of telomeres (ALT) is a major determinant of poor outcome in neuroblastoma. Here, we screen for ALT in primary and [Read more](#) →

Study 1 / 5

[Next Study](#)

Published in:



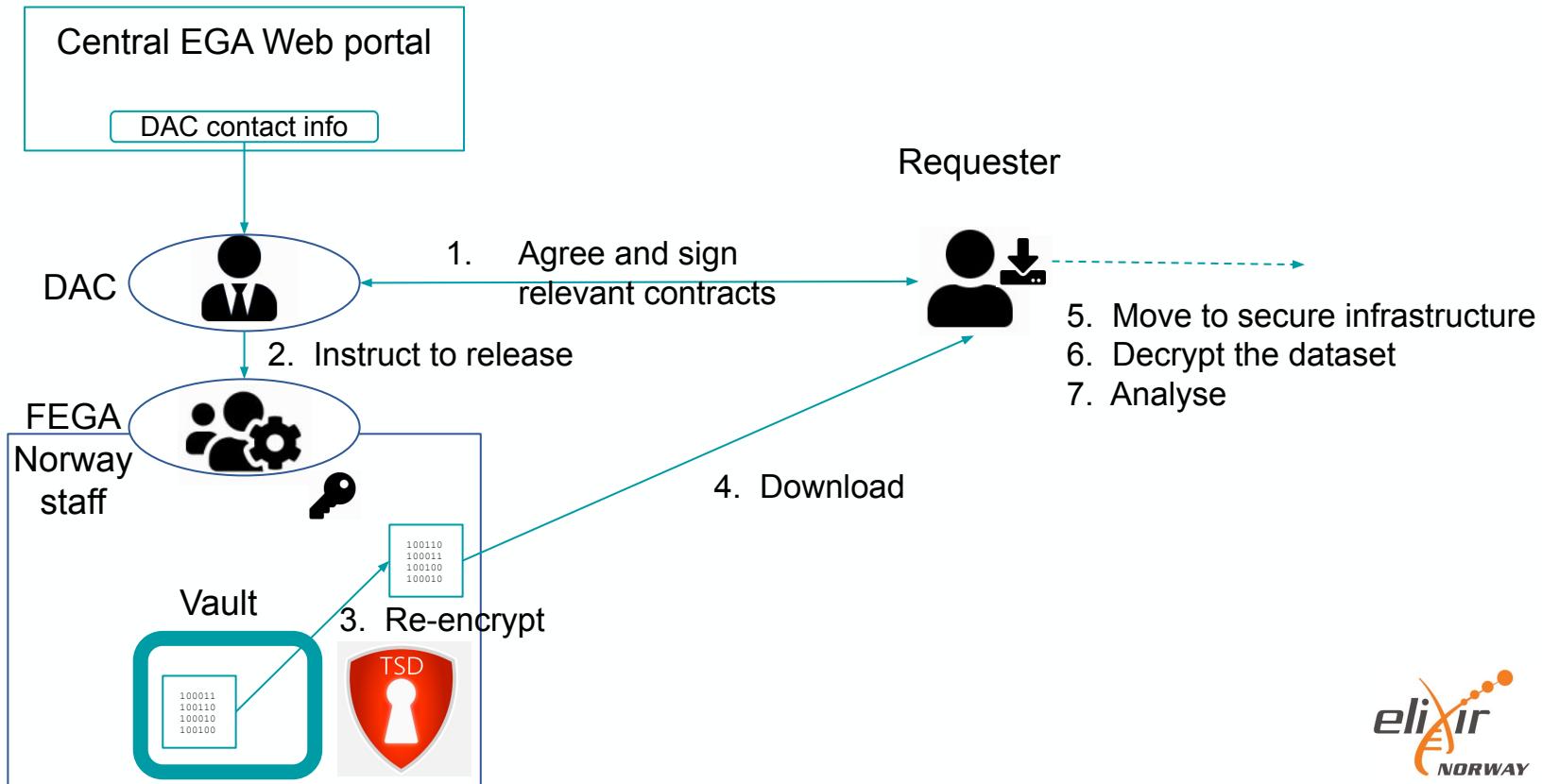
BROWSE

[Studies](#)[Datasets](#)[DACs](#)

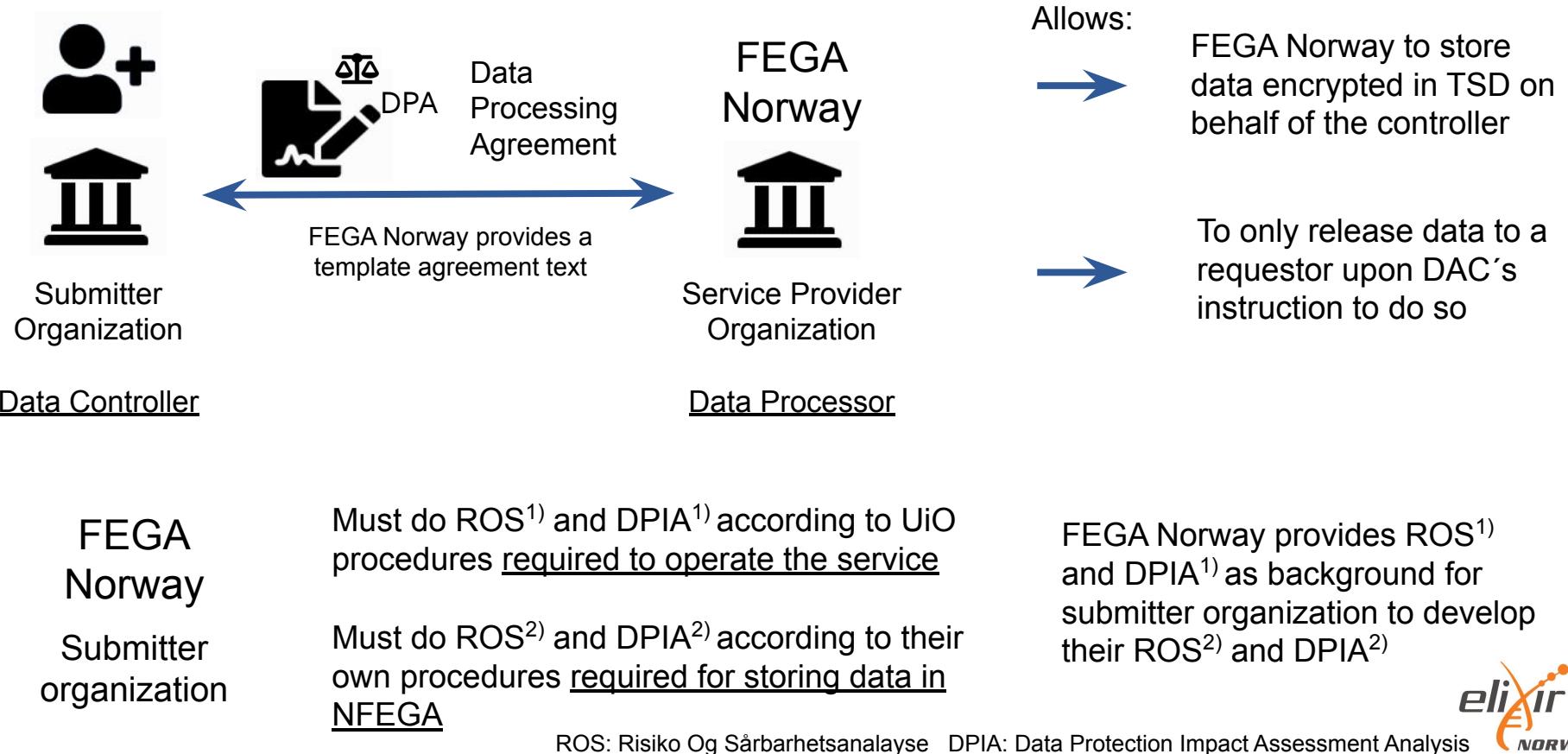
HELP

[FTP & Aspera](#)[Tools](#)[EGA Blog](#)

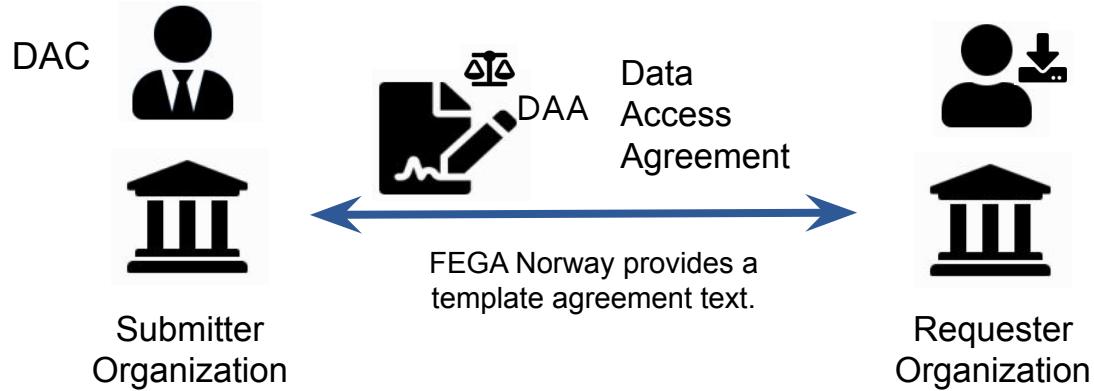
Data flow and data control - part II



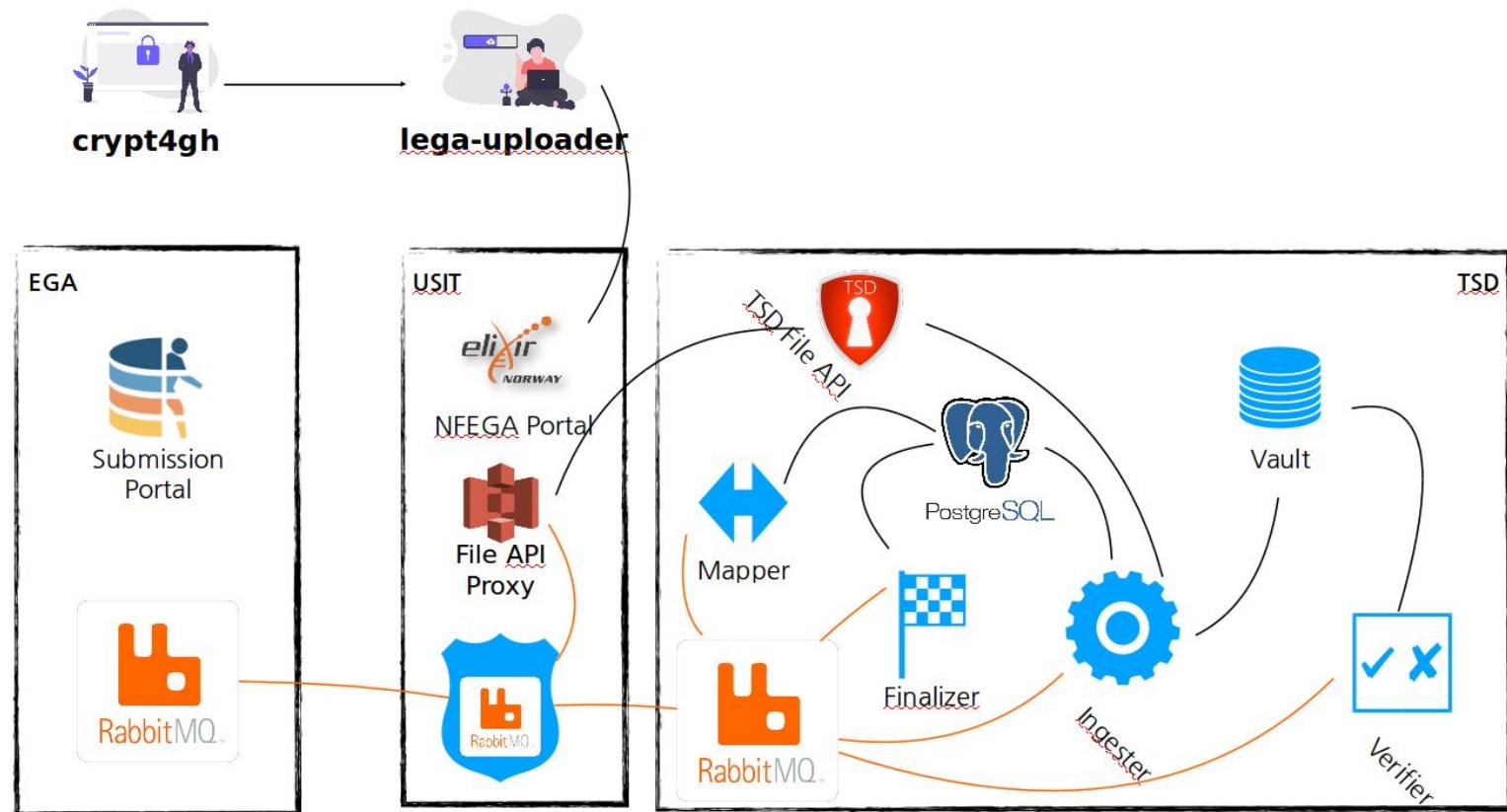
GDPR required agreements and assessments



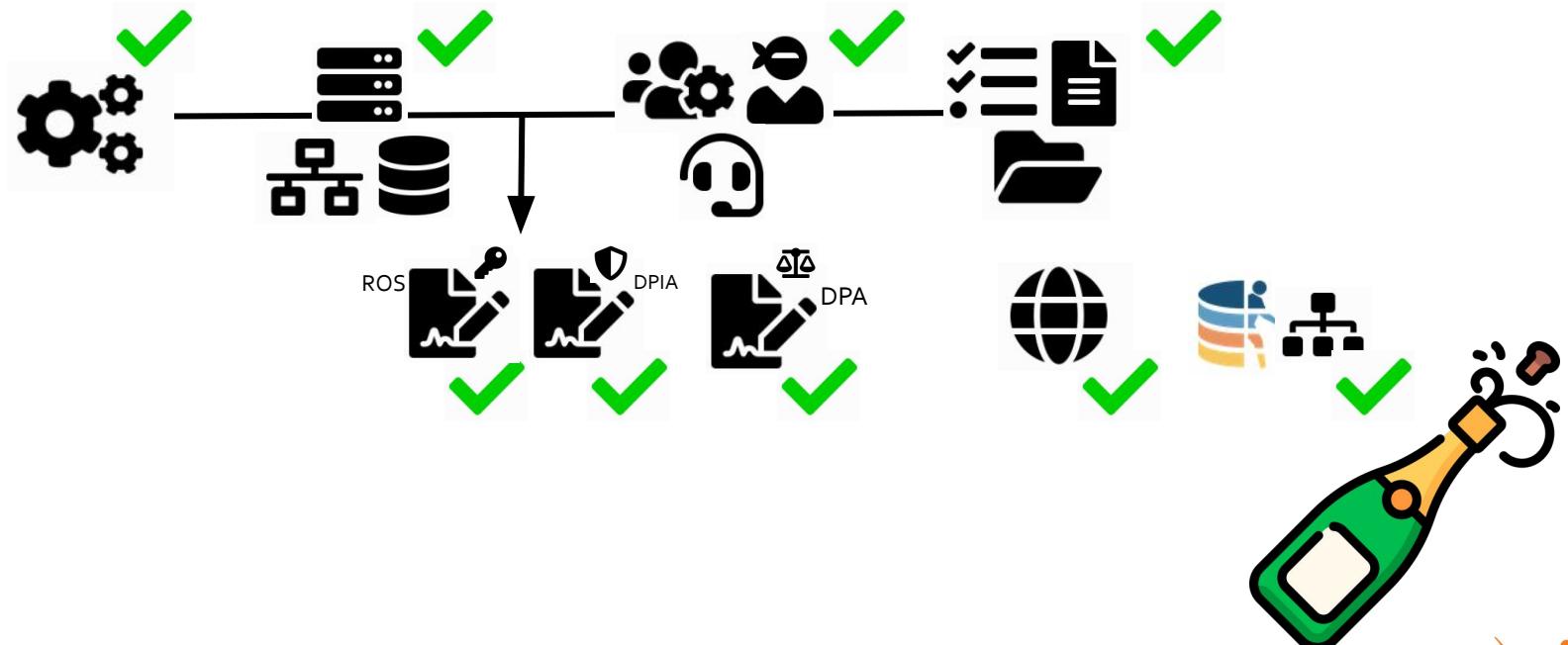
Data Access Agreement before releasing the data to Requester



FEGA Norway- Data-In



What is already in place for FEGA Norway ?



23.05.2022

First submitter with dataset ready to be uploaded



Janus Serum Bank

The Janus Serum Bank has blood samples from 318 628 Norwegians. The biobank is reserved for cancer research, and is internationally unique regarding size and number of cancer cases.



The logo for the Janus Serum Bank is circular. Inside the circle, there is a stylized illustration of a human head facing left, with a cluster of small circles representing a brain or genetic material. The text "Janus Serum Bank" is written in a blue, sans-serif font along the top inner edge of the circle.

- “Small non-coding RNA as early detector cancer biomarkers”
- Close to 2000 Samples
 - 10 different cancer types
 - 1 control group
- miRNA-seq data
- DAC appointed

How can XUH or UiX support their researchers with FEGA data submission?

- DPA (+DPIA and ROS) of “Submitting my research dataset to FEGA Norway” is daunting for individual researchers
 - Involves competence from and communication with many OUS units/roles:
 - Data Protection Officer (Personvernombud)
 - Data Security Officer
 - Legal advisors
- DPA and related work is specific to the dataset in question, but there are large parts that are common for all datasets submitted to FEGA Norway.
 - Can be handled efficiently at right organizational level through templates and standardised agreements with appendix for dataset specific details.



Protecting

The sensitive data stored at Federated EGA instance is encrypted with Crypt4GH - a new standard file container format from the Global Alliance for Genomics and Health (GA4GH), allowing genomic data to remain secure throughout their lifetime, from initial sequencing to sharing with professionals at external organizations.

Strong encryption both in transit and at rest are critical requirements for Genomics England and other genomics in healthcare initiatives.

Augusto Rendon, Crypt4GH implementation leader at Genomics England

Hosted by



Deposit data

Please follow the instructions on how to perform a [submission](#).



Search data

Please follow the instructions on how to perform a [search](#).



Retrieve data

Please follow the instructions on how to perform a [retrieval](#).



Approve access

DAC, please follow the instructions of data [approval](#).

ega-norway-support@elixir.no

Design by [Papaya](#). Illustrations from [Undraw](#).



Where to work on your data before publication?



NTNU | HUNT Cloud

Acknowledgements



Eivind Hovig
Fabian L. M. Bernal
Milen Kouylekov
Ahmed Ghanem
Arash Azamifard



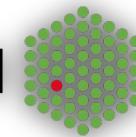
Gard Thomassen
Leon Charl du Toit
Abdulrahman Azab



HEILSA TRYGGVEDOTTIR



EMBL-EBI



Dmytro Titov
(former team member)





Useful Resources

- Federated EGA Norway node, <https://ega.elixir.no>
- TSD (UiO) - <https://www.uio.no/tsd>
- SAFE (UiB) - <https://wwwuibno/safe> (English link top right)
- Hunt Cloud (NTNU) - <https://www.ntnu.edu/mh/huntcloud>
- RDMkit Sharing step: <https://rdmkit.elixir-europe.org/sharing>

Licensing



Learning Objectives

In this part, we will learn about data sharing with respect to:

- code versioning
- code deposition
- licensing
- data availability

Instructor: Korbinian Bösl



Learning Activity

Data Sharing [Quiz Link](#)

Data life cycle	+
Your role	+
Your domain	+
Your problem	-

Compliance monitoring

Data analysis

Data management plan

Data organisation

Data protection

Data publication

Data quality

Data storage

Data transfer

Identifiers

Licensing

Documentation and metadata

Sensitive data

All tools and resources

Tool assembly



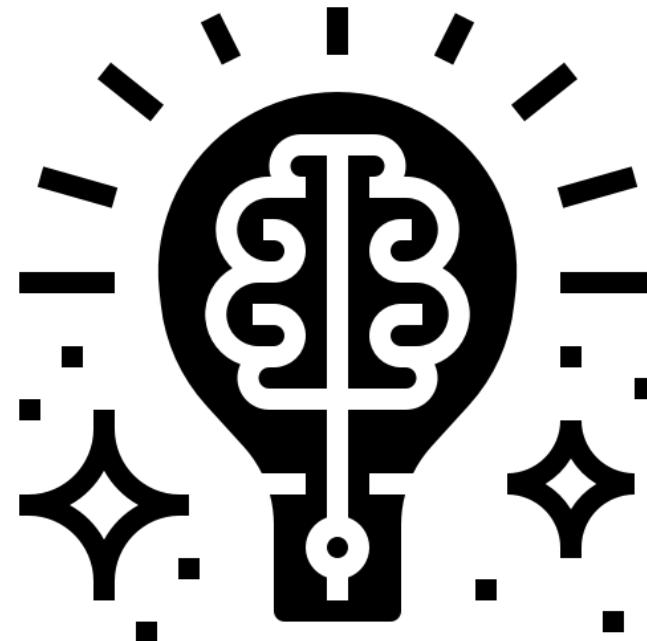
Link to RDMkit:
<https://rdmkit.elixir-europe.org/>

Intellectual property rights (IPRs)

Patent: Protects novel, non-obvious, inventions

Copyright: creative products: software, writing, figures, photos, some datasets, this presentation

Trademark: Protects a name/brand



CC-BY 3.0 Wichai.Wi

IPR often regulated in work contract

Why should I licence my research outputs?

Why should I licence my research outputs?



Legal security for users (Accessibility)

Why should I licence my research outputs?



Legal security for users (Accessibility)



Increase of willingness to reuse outputs (Reusability)

Why should I licence my research outputs?



Legal security for users (Accessibility)



Increase of willingness to reuse outputs (Reusability)



Allows deposition/mirroring in 2nd databases (Findability)

Concepts in open licenses



Waive all your interests that may exist in your work

Copy left:

Concepts in open licenses



Waive all your interests that may exist in your work

Copy left:



Credit for the original creation

Concepts in open licenses



Waive all your interests that may exist in your work

Copy left:



Credit for the original creation



License new creations under identical terms

Concepts in open licenses



Waive all your interests that may exist in your work

Copy left:



Credit for the original creation



License new creations under identical terms



Non-commercial

Concepts in open licenses



Waive all your interests that may exist in your work

Copy left:



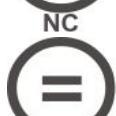
Credit for the original creation



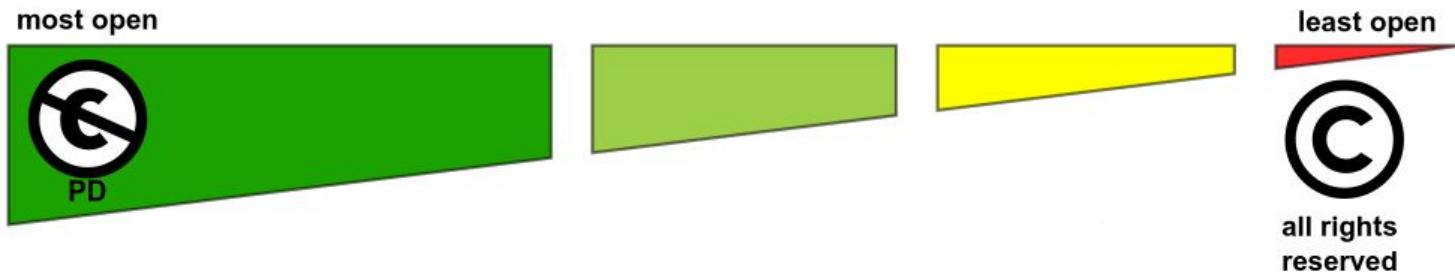
License new creations under identical terms

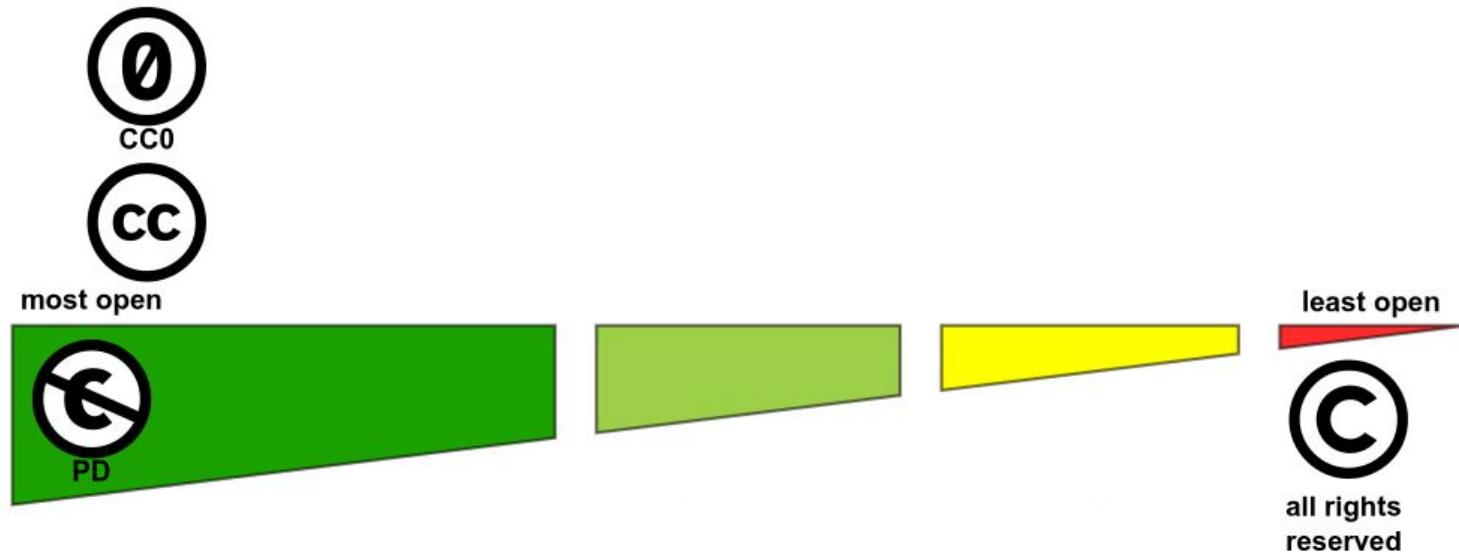


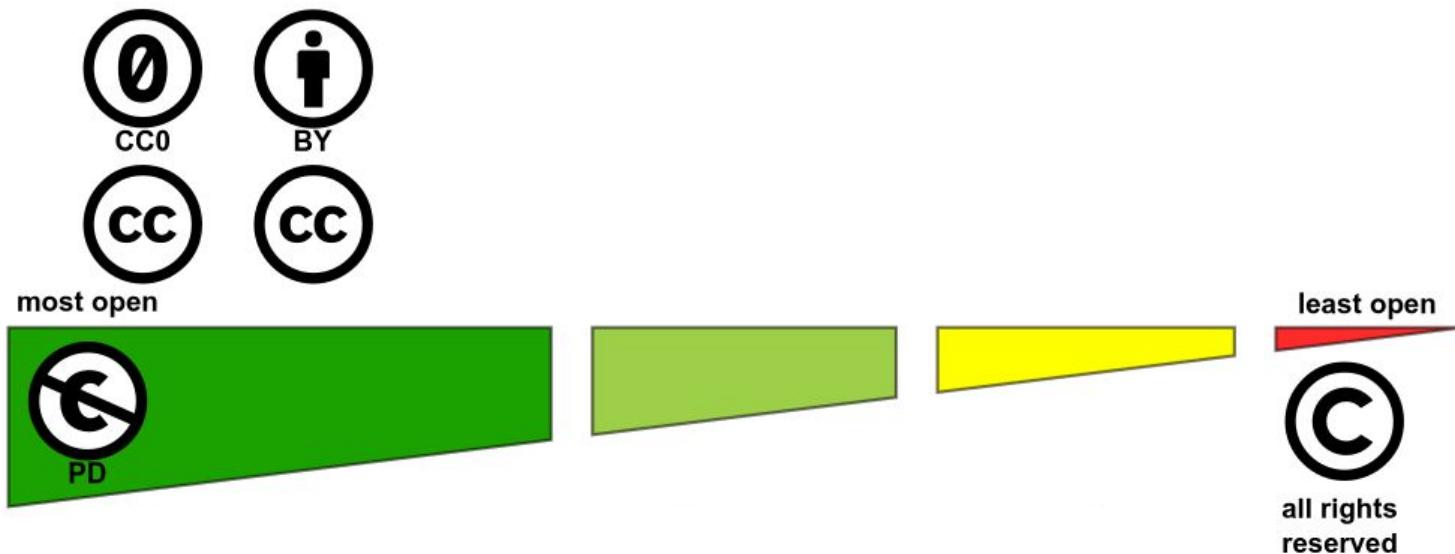
Non-commercial

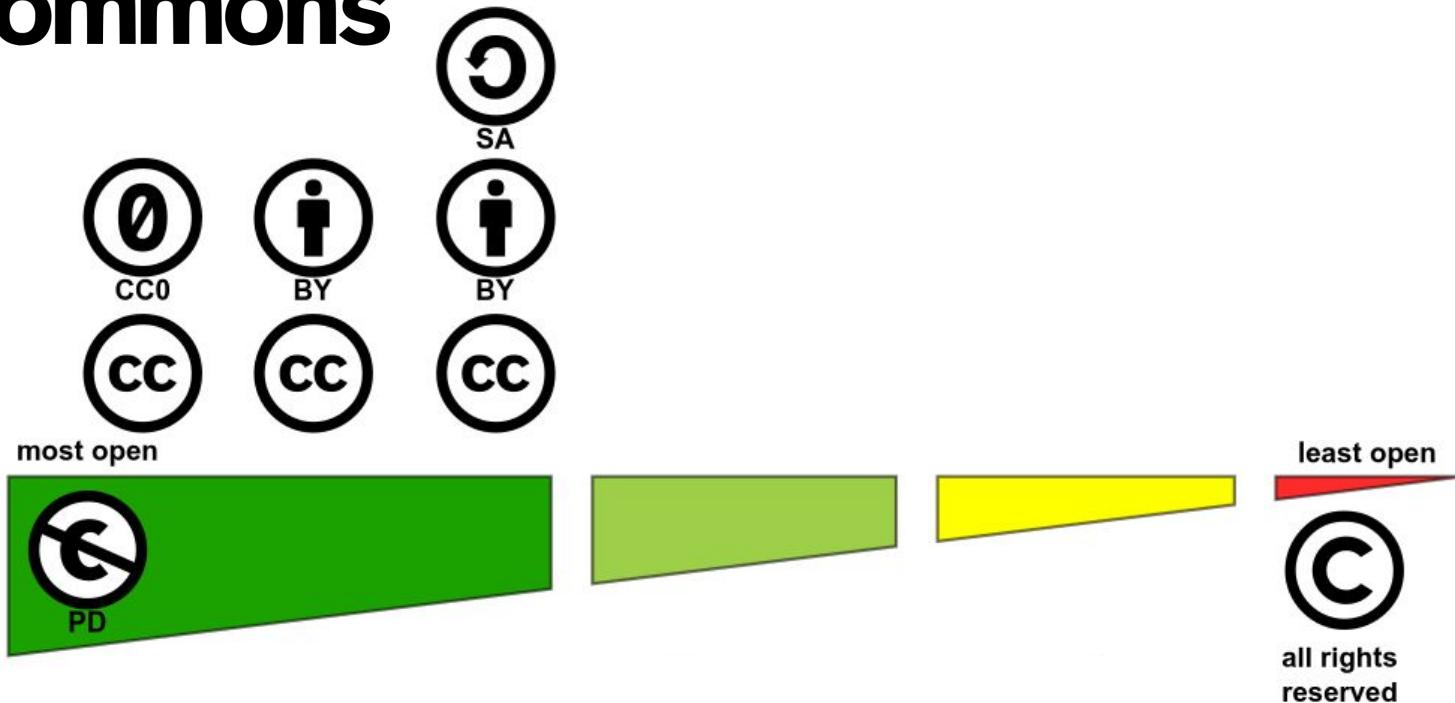


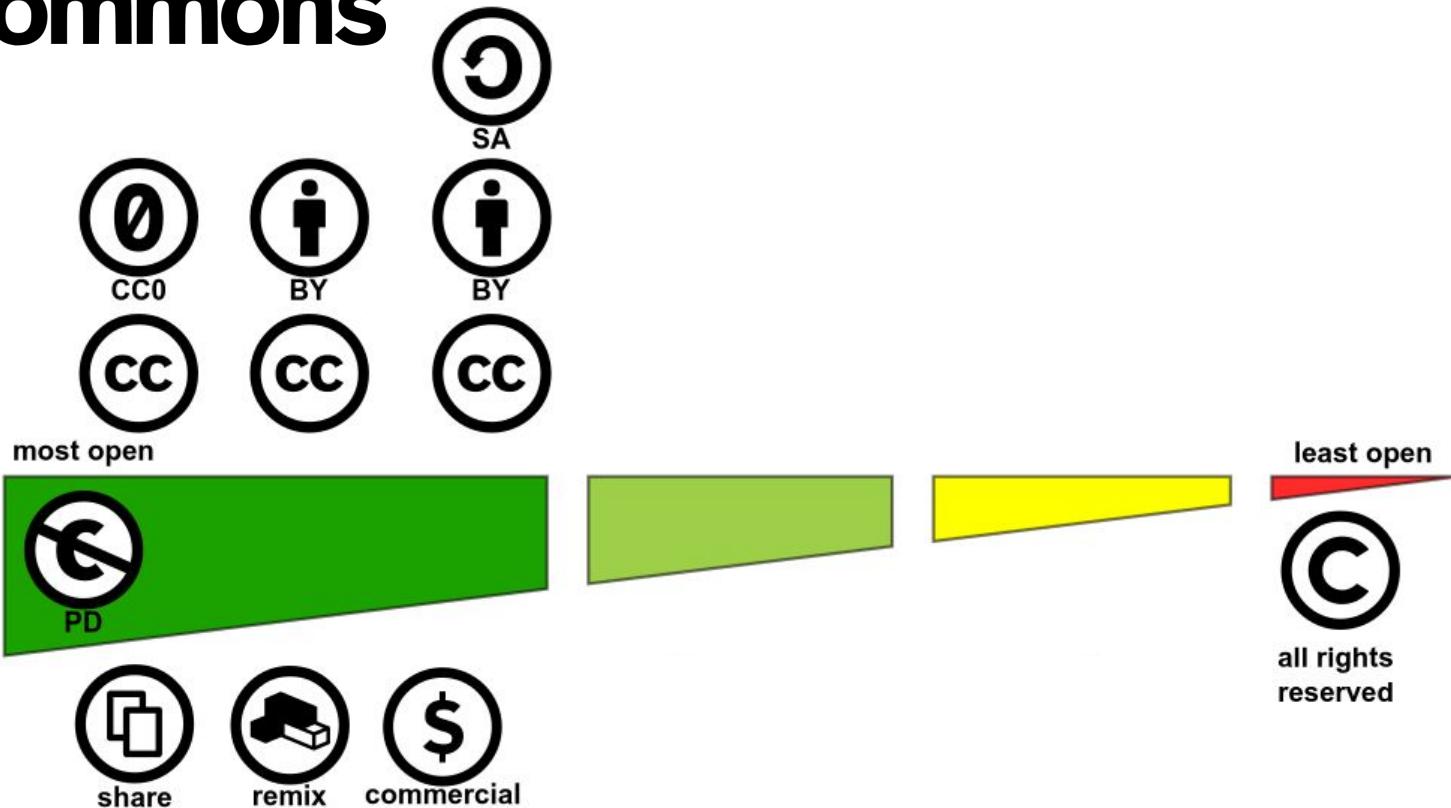
Cannot be shared with others in adapted form

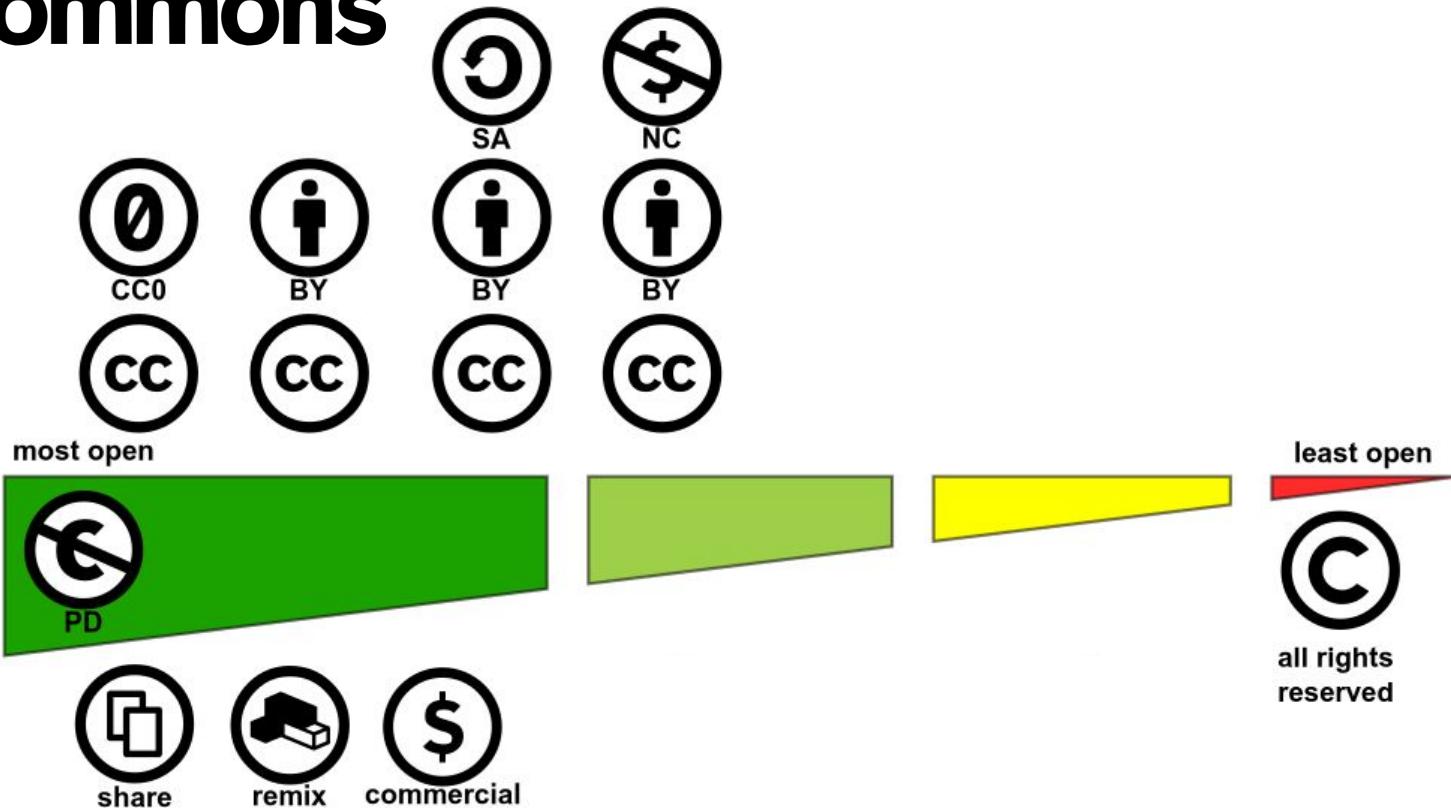


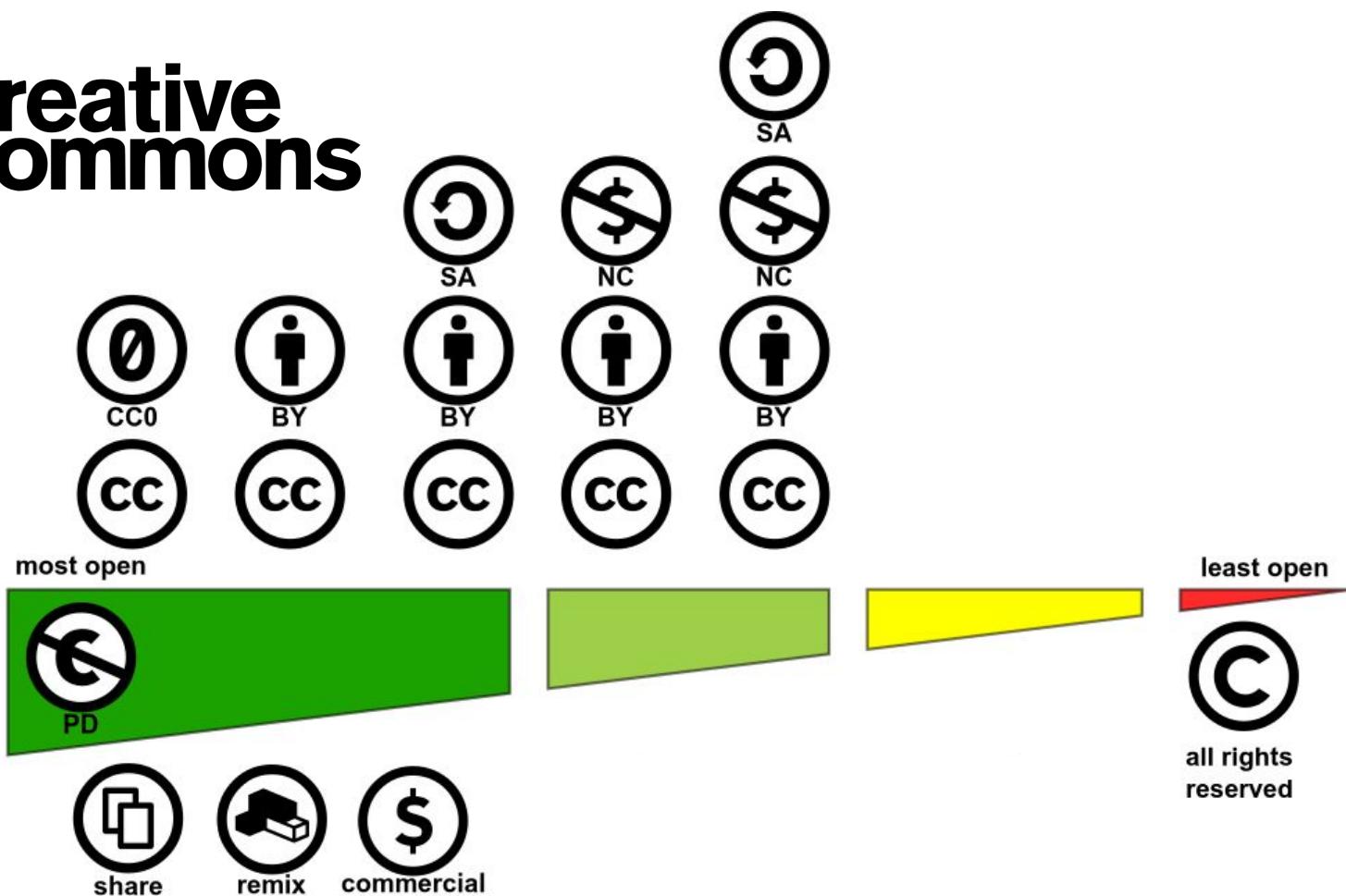


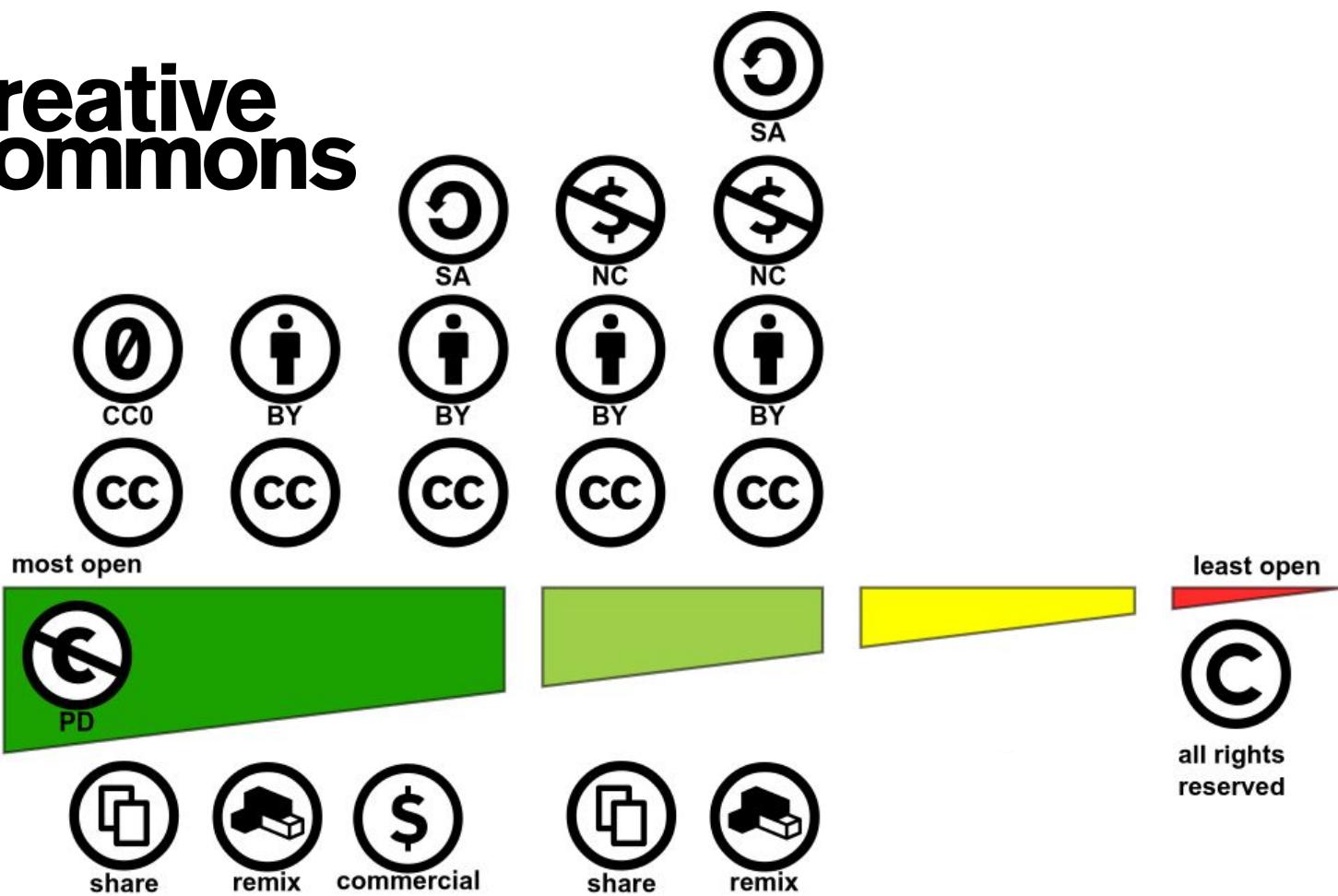


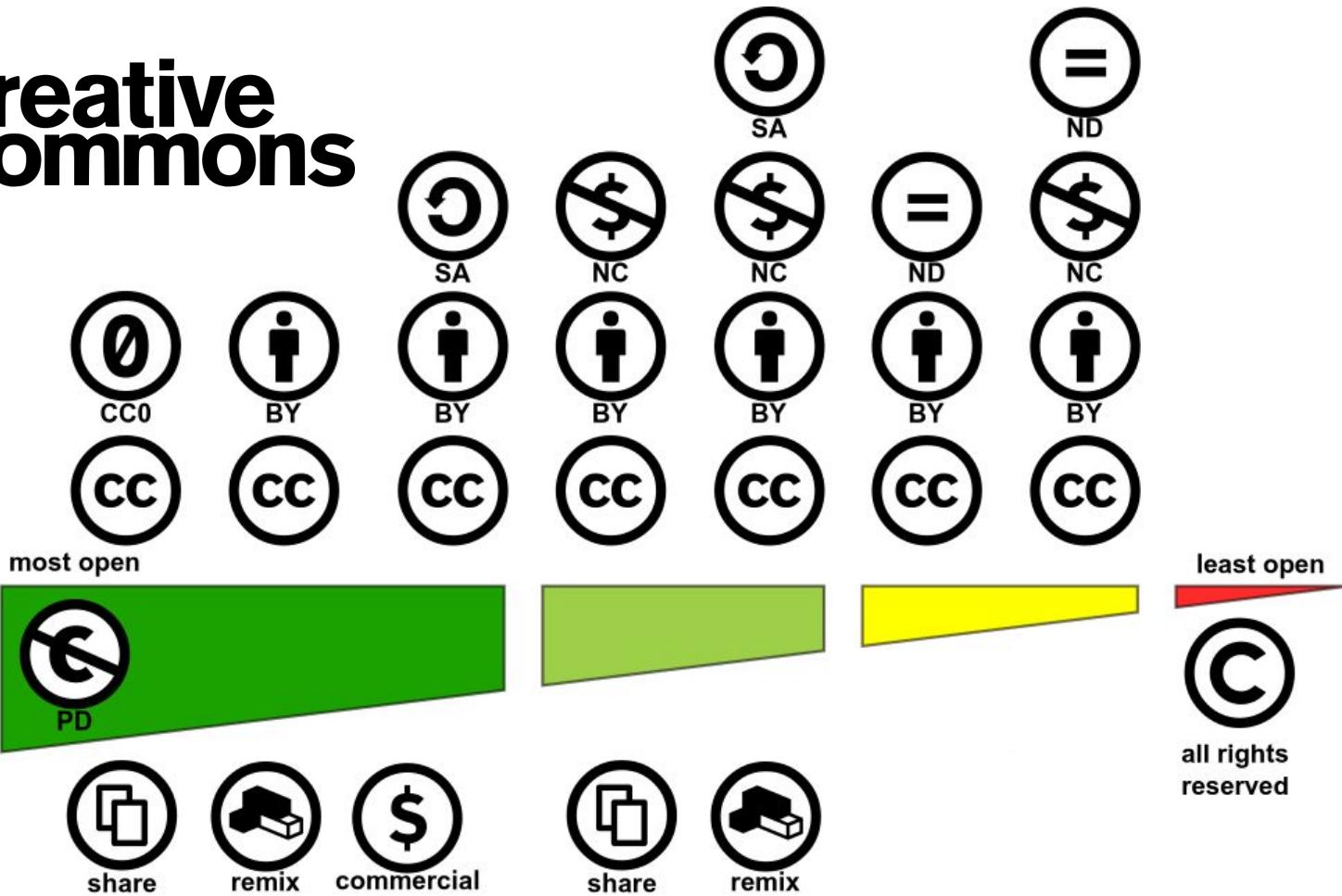


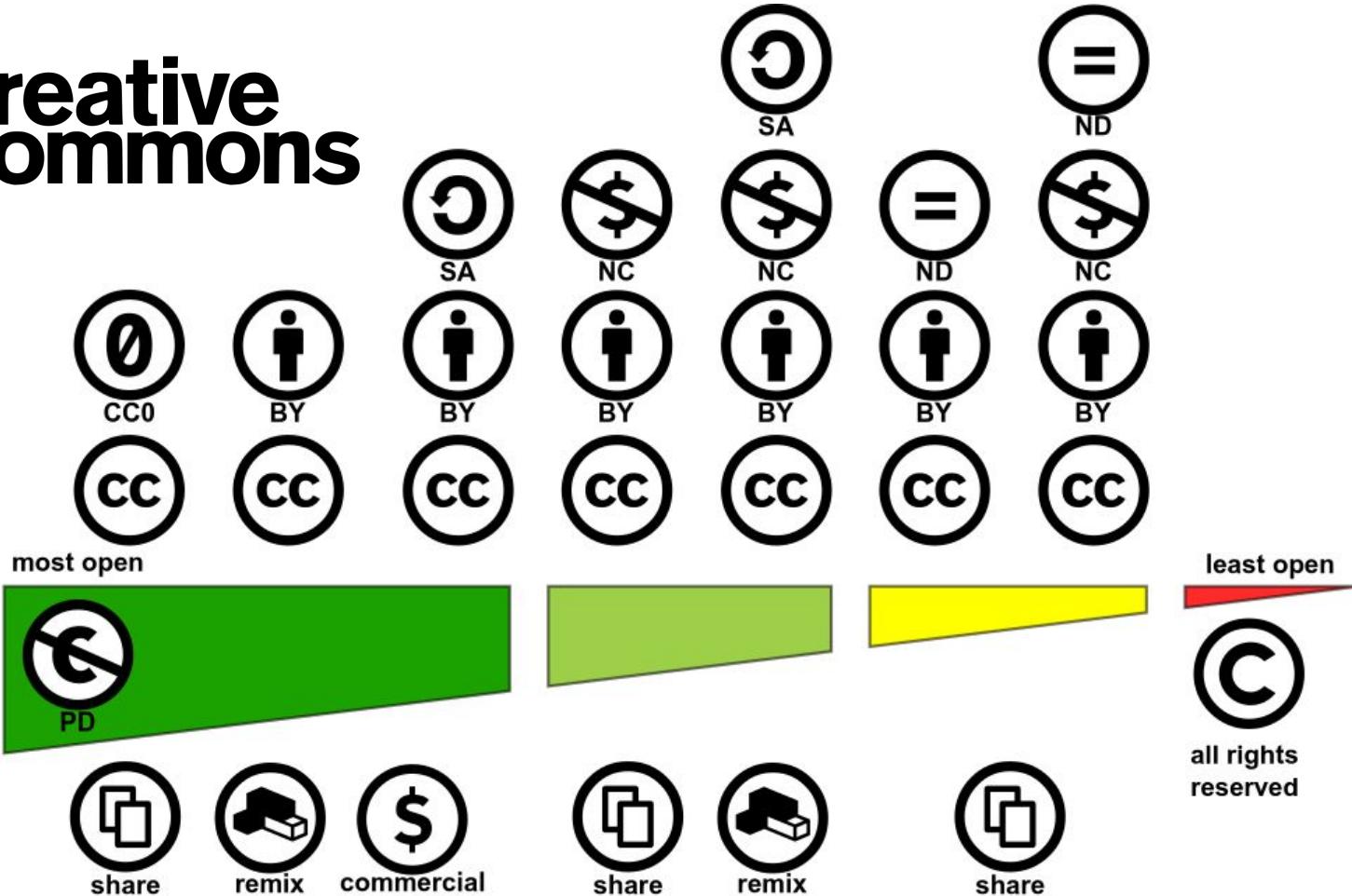




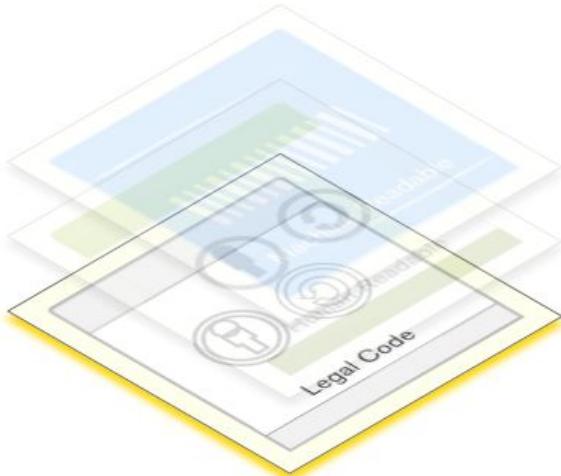












Legal code, harmonized for
national/international law



Legal code, harmonized for
national/international law

Creative Commons - **What is licensed?** readable, understandable text



Nathan Yergler, Alex Roberts - **Who is to be attributed?**

Licensed to the public under [CC BY 3.0 Unported](#) - **Which license?**



**Legal code, harmonized for
national/international law**

Human readable, understandable text

**Machine readable html tag attachable
to metadata**



Challenges



Multiple Attributions for several sources (license stacking)



Challenges



Challenges



Multiple Attributions for several sources (license stacking)

BY



Multiple incompatible source licenses

SA



Challenges



Multiple Attributions for several sources (license stacking)



Multiple incompatible source licenses



Legal commercial definition (e.g. use by journals)



Challenges



Multiple Attributions for several sources (license stacking)



Multiple incompatible source licenses



Legal commercial definition (e.g. use by journals)



Unintentional restrictive

Open Data Commons

Databases are different to simple data (e.g. EU-copyright)

Open Data Commons

Databases are different to simple data (e.g. EU-copyright)



Open Data Commons Public Domain Dedication and License
(PDDL)

Open Data Commons

Databases are different to simple data (e.g. EU-copyright)



**Open Data Commons Public Domain Dedication and License
(PDDL)**



Open Data Commons Attribution License (ODC-By)

Open Data Commons

Databases are different to simple data (e.g. EU-copyright)



Open Data Commons Public Domain Dedication and License
(PDDL)



Open Data Commons Attribution License (ODC-By)



Open Data Commons Open Database License (ODbL)

Repository specific regulations



Individuals submitting data to the international sequence databases managed collaboratively by DDBJ, EMBL, and GenBank should be aware of the following:

Repository specific regulations



Individuals submitting data to the international sequence databases managed collaboratively by DDBJ, EMBL, and GenBank should be aware of the following:

The INSDC has a uniform policy of **free and unrestricted access** to all of the data records their databases contain. Scientists worldwide can access these records to plan experiments or publish any analysis or critique. Appropriate credit is given by citing the original submission, following the practices of scientists utilising published scientific literature.

The INSDC will **not attach statements to records that restrict access to the data**, limit the use of the information in these records, or prohibit certain types of publications based on these records. Specifically, no use restrictions or licensing requirements will be included in any sequence data records, and no restrictions or licensing fees will be placed on the redistribution or use of the database by any party.

Norwegian Licence for Open Government Data (NLOD) 2.0



Digitaliseringsdirektoratet
data.norge.no

Norwegian Licence for Open Government Data (NLOD) 2.0



Digitaliseringsdirektoratet
data.norge.no

A licence compatible by contract shall mean the following licences:

for all information: Open Government Licence (version 1.0, 2.0 and 3.0), **Creative Commons Attribution Licence (international version 4.0 and norwegian version 4.0)**

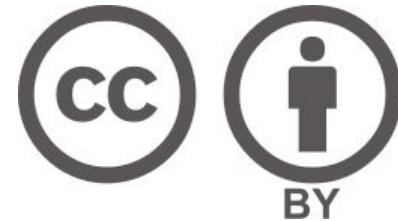
for those parts of the information which do not constitute databases: **Creative Commons Attribution Licence (generic version 1.0, 2.0, 2.5 and unported version 3.0)** and **Creative Commons Navngivelse 3.0 Norge**

for those parts of the information which constitute databases: **Open Data Commons Attribution License (version 1.0)**.

Norwegian Licence for Open Government Data (NLOD) 2.0



Digitaliseringsdirektoratet
data.norge.no



A licence compatible by contract shall mean the following licences:

for all information: Open Government Licence (version 1.0, 2.0 and 3.0), **Creative Commons Attribution Licence (international version 4.0 and norwegian version 4.0)**

for those parts of the information which do not constitute databases: **Creative Commons Attribution Licence (generic version 1.0, 2.0, 2.5 and unported version 3.0)** and **Creative Commons Navngivelse 3.0 Norge**

for those parts of the information which constitute databases: **Open Data Commons Attribution License (version 1.0)**.

Open Source Software licenses

Special considerations for Software

Liability

Warranty

Modifications

Network use = Distribution?



Open Source Software licenses





Open Source Software licenses



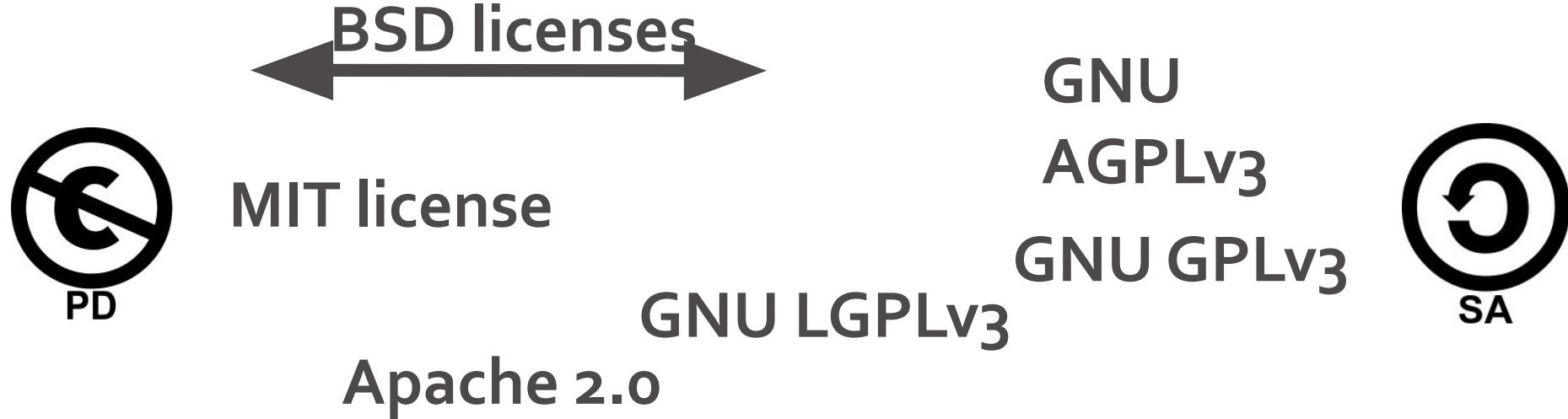
MIT license

GNU
AGPLv3



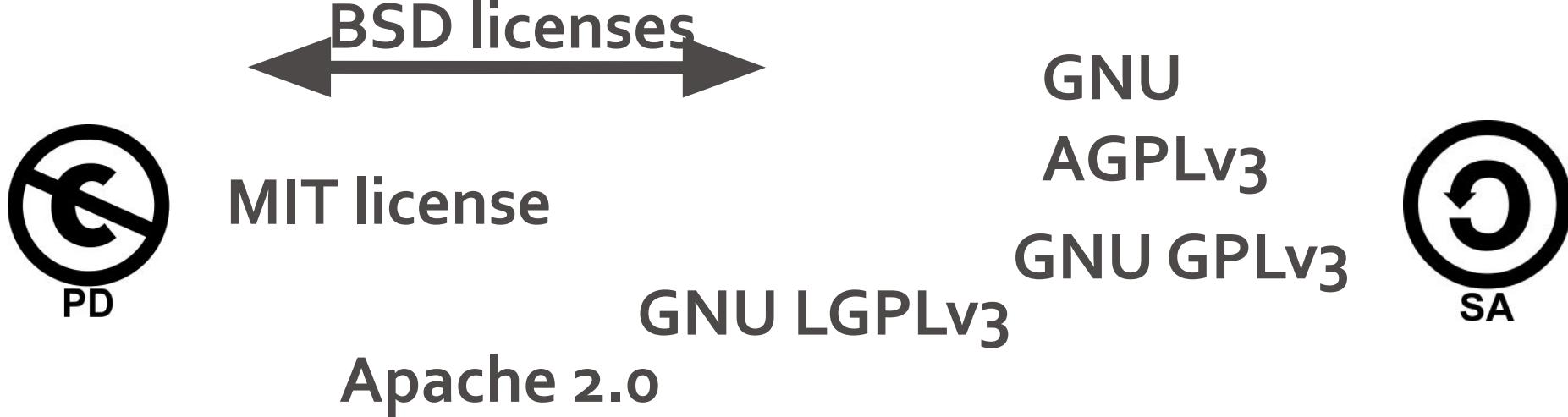


Open Source Software licenses





Open Source Software licenses



<https://opensource.org/licenses>

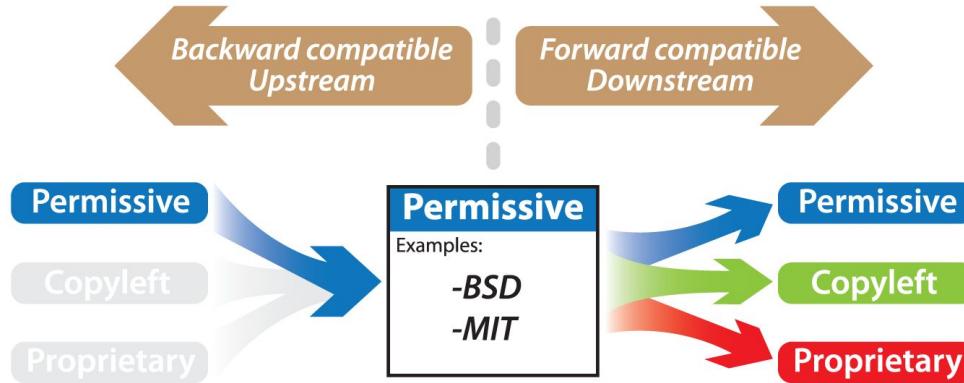
<https://choosealicense.com/>



Icons CC-BY 4.0
https://creativecommons.org/

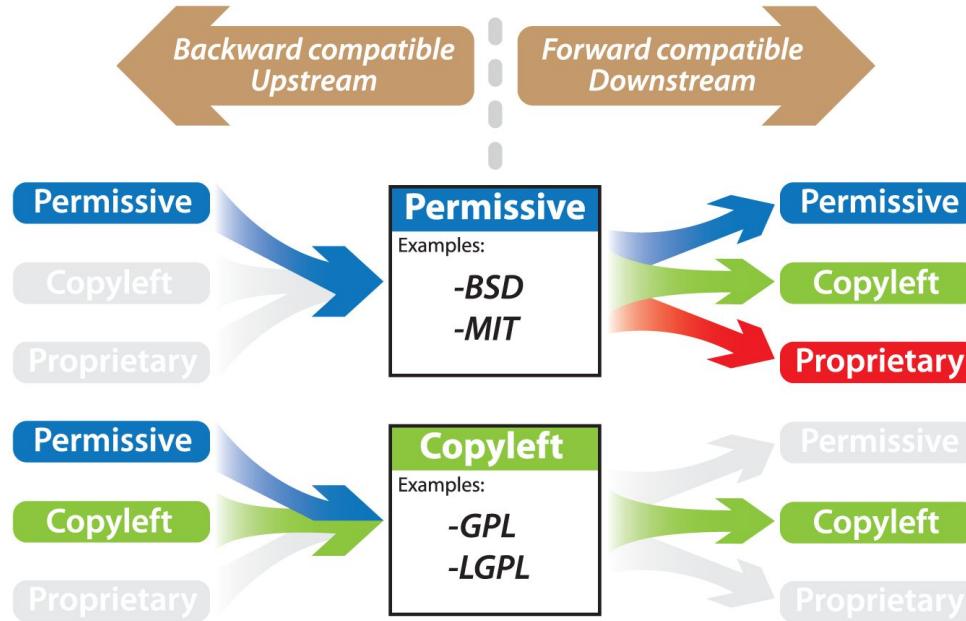


Open Source Software licenses



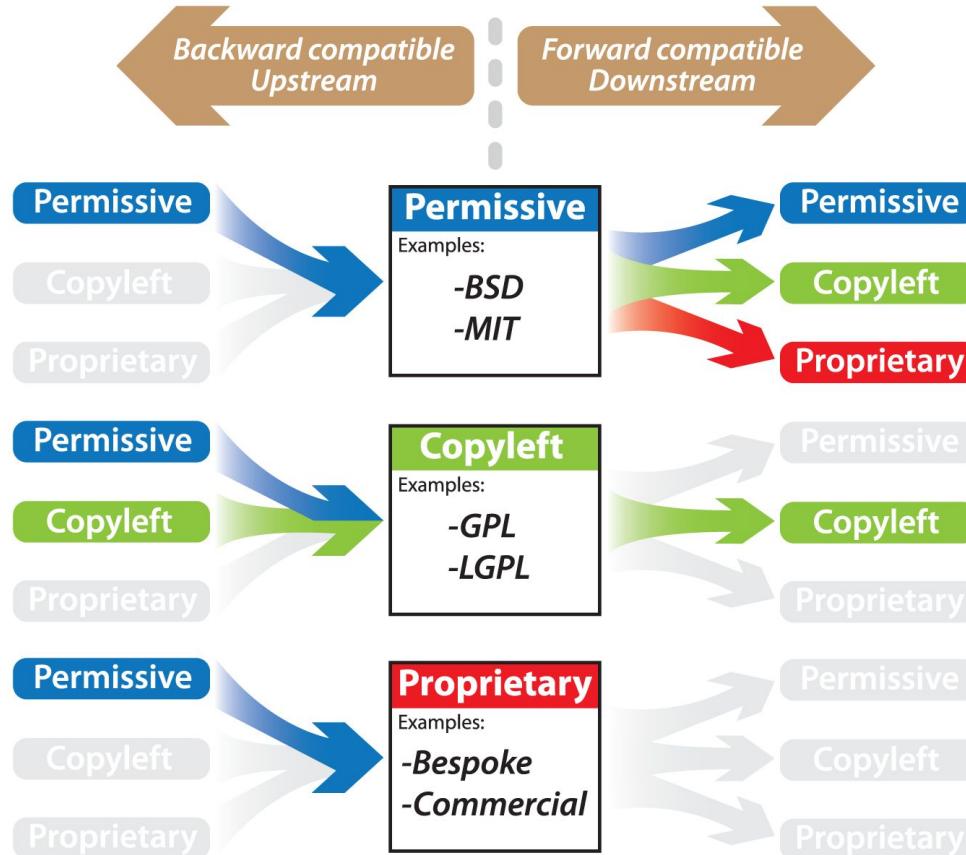


Open Source Software licenses





Open Source Software licenses



Material Transfer Agreements

Used for e.g. reagents, cell lines, plasmids, mice, ...

Can safeguard your commercial interest, while allowing others to use your material for research

Can ensure attribution

Can enforce remain in Public Domain

Uniform Biological Material Transfer Agreement



Benefits to the Recipient:

Permission to use material for research or teaching purposes

Rights to all research results, modifications, and invention

Patent applications on modifications or inventions

Publishing without editorial comment or review by provider

Limited liability

Data availability statement

“Data available upon
reasonable request”

... so it's not?

Gabelica, M., Bojčić, R. & Puljak, L. *Journal of Clinical Epidemiology* **0**, (2022). <https://doi.org/10.1016/j.jclinepi.2022.05.019>

Tedersoo, L. et al. *Sci Data* **8**, 192 (2021). <https://doi.org/10.1038/s41597-021-00981-0>

Data availability statement

RNA-Seq data are available in

the ArrayExpress database

<http://www.ebi.ac.uk/arrayexpress/experiments/E-MEXP-12345/>

under accession number E-MEXP-12345

<https://www.ebi.ac.uk/arrayexpress/about.html>

Data availability statement

RNA-Seq data are available in - which dataset?

the ArrayExpress database - which repository?

<http://www.ebi.ac.uk/arrayexpress/experiments/E-MEXP-12345/> - Link!

under accession number E-MEXP-12345 - ID

<https://www.ebi.ac.uk/arrayexpress/about.html>

Data availability statement

RNA-Seq data are available in - which dataset?

the ArrayExpress database - which repository?

<http://www.ebi.ac.uk/arrayexpress/experiments/E-MEXP-12345/> - Link!

under accession number E-MEXP-12345 - ID

(+access conditions & license)

-> Check repository and journal guidelines!

<https://www.ebi.ac.uk/arrayexpress/about.html>



Useful Resources

- Federated EGA Norway node, <https://ega.elixir.no>
- TSD (UiO) - <https://www.uio.no/tsd>
- SAFE (UiB) - <https://www.uib.no/safe> (English link top right)
- Hunt Cloud (NTNU) - <https://www.ntnu.edu/mh/huntcloud>
- Data Sharing <https://rdmkit.elixir-europe.org/sharing>
- Data Licensing <https://rdmkit.elixir-europe.org/licensing>
- Data Brokering https://rdmkit.elixir-europe.org/data_brokering
- [EUDAT license chooser for data and software](#)



Learning Activity

Data Sharing [Quiz Link](#)

Thank you!



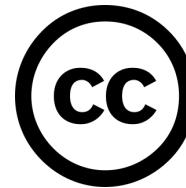
elixir-norway.org



@elixirnorway



contact@bioinfo.no



Except where otherwise noted, this work is licensed under a

Creative Commons Attribution 4.0 International License

<https://creativecommons.org/licenses/by/4.0/>