

# Data brokering

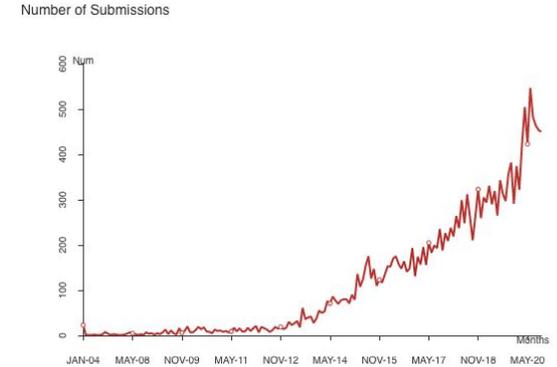
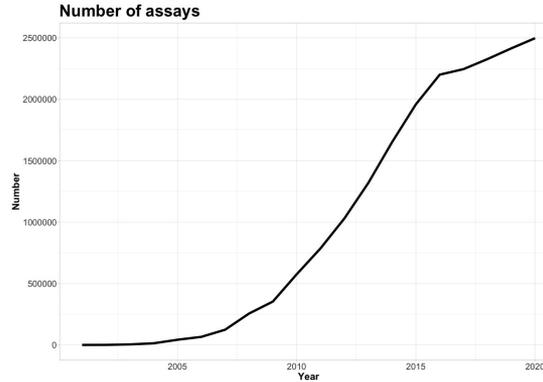
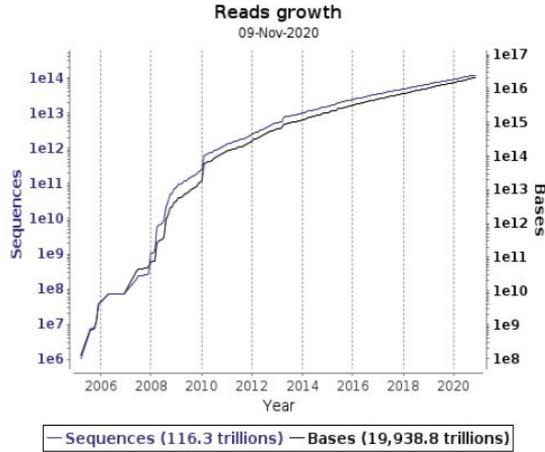
[frama.link/ifb-ag20-data-brokering](https://frama.link/ifb-ag20-data-brokering)

Thomas Denecker, H el ene Chiapello & Jacques van Helden

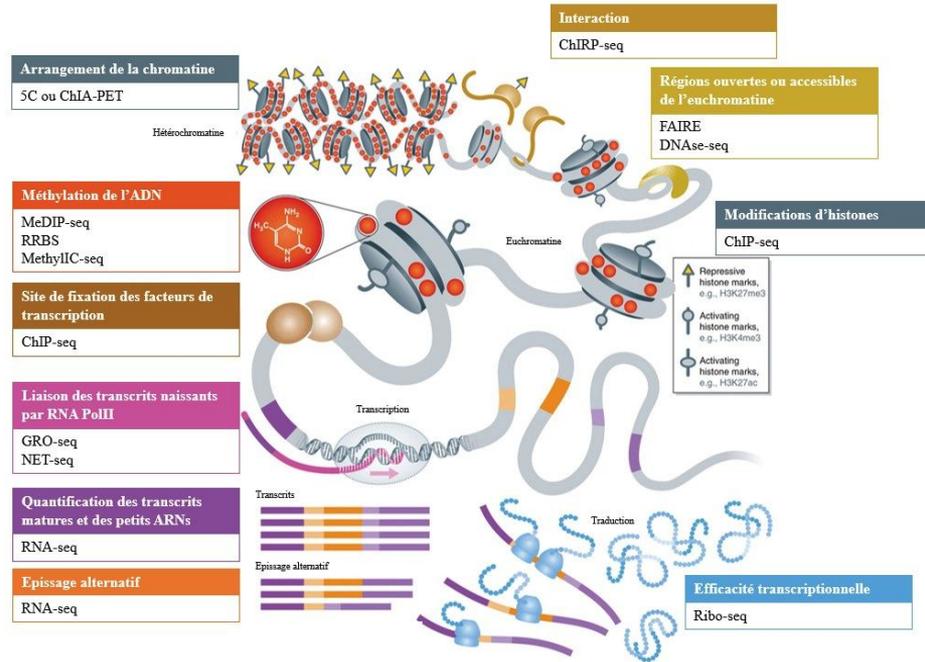


*Contexte*

## Toujours plus de données...



... de plus en plus hétérogènes ...



(Soon et al. 2013)

... stockées dans des bases de données de plus en plus nombreuses



≈ **1700** bases de données  
référéncées dans NAR Database en 2018

(Rigden et al, 2018)

*Soumission*

## Soumission

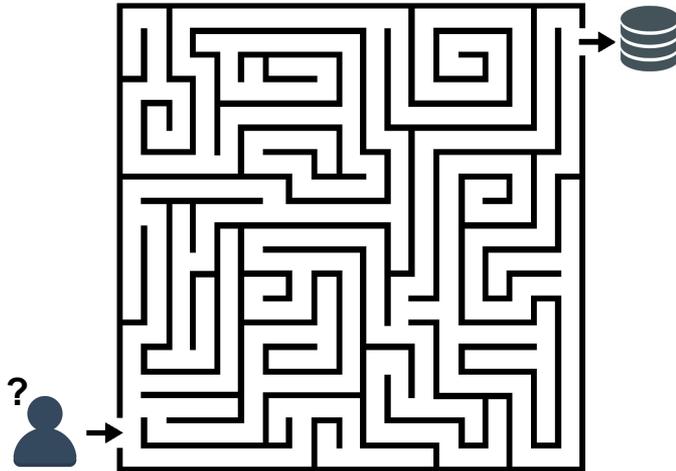


### Aujourd'hui

Soumission par l'équipe  
expérimentale ou bioinformatique

Allers-retours entre la base de  
données et l'utilisateur

## Soumission du point de vue de l'utilisateur



**Le sentiment général ?**

*“Pas toujours si simple”*

**Des questions récurrentes**

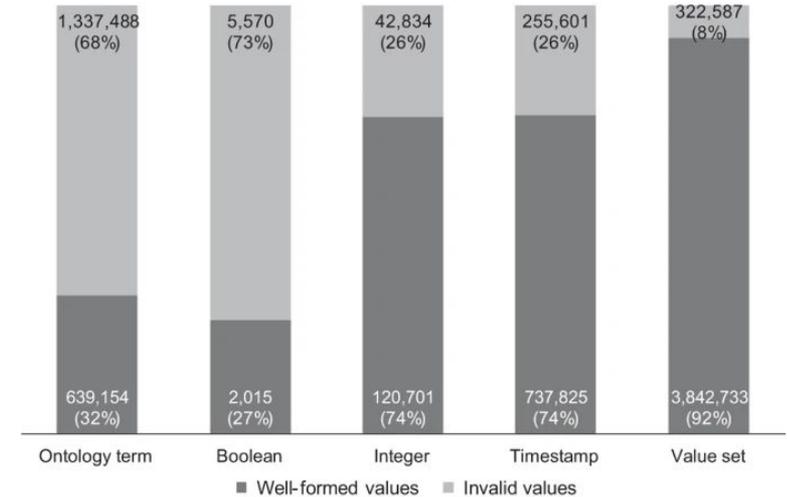
*“Dans quelle base de données soumettre ?  
Quelles données soumettre ? Les brutes, les  
traitées ? Accompagnées de quelles  
informations supplémentaires ? ...”*

## Bilan après soumission

### Manque de consistance des soumissions

### Métadonnées de mauvaise qualité

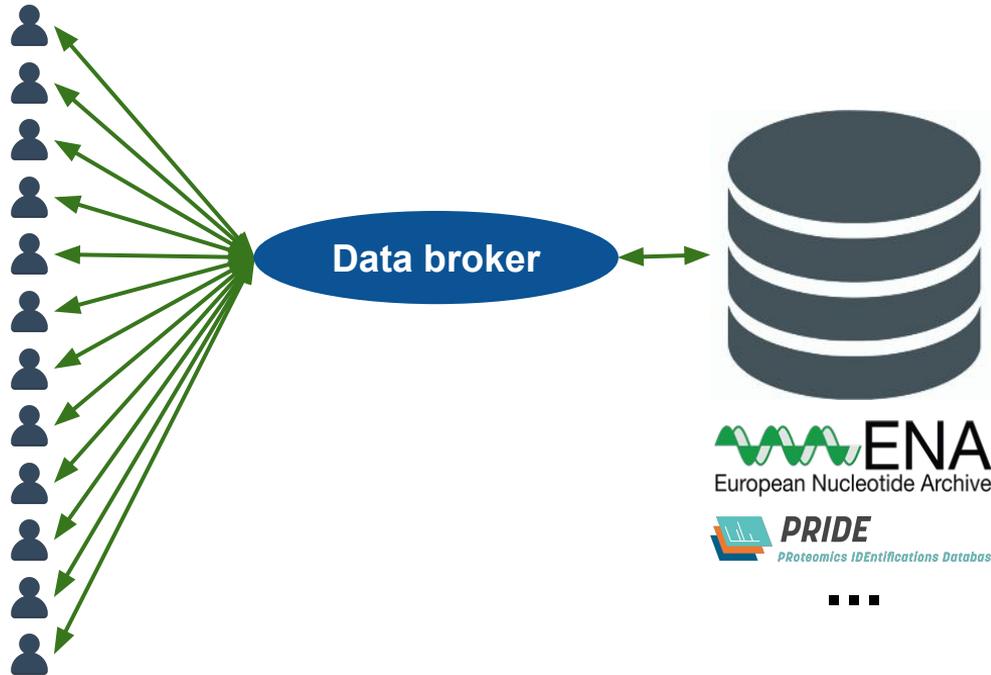
- Incohérentes
- Incomplètes
- Peu informatives
- Redondantes



(Gonçalves et al., 2019)

# *Data brokering*

## Un data broker



Intermédiaire entre le  
producteur de données et la  
ressource internationale de  
stockage/archivage

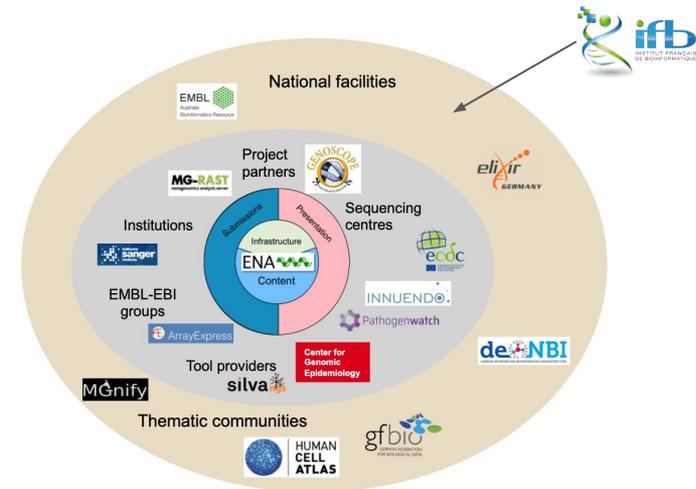
Rationalisation des échanges

## Data broker & Soumission

### Missions

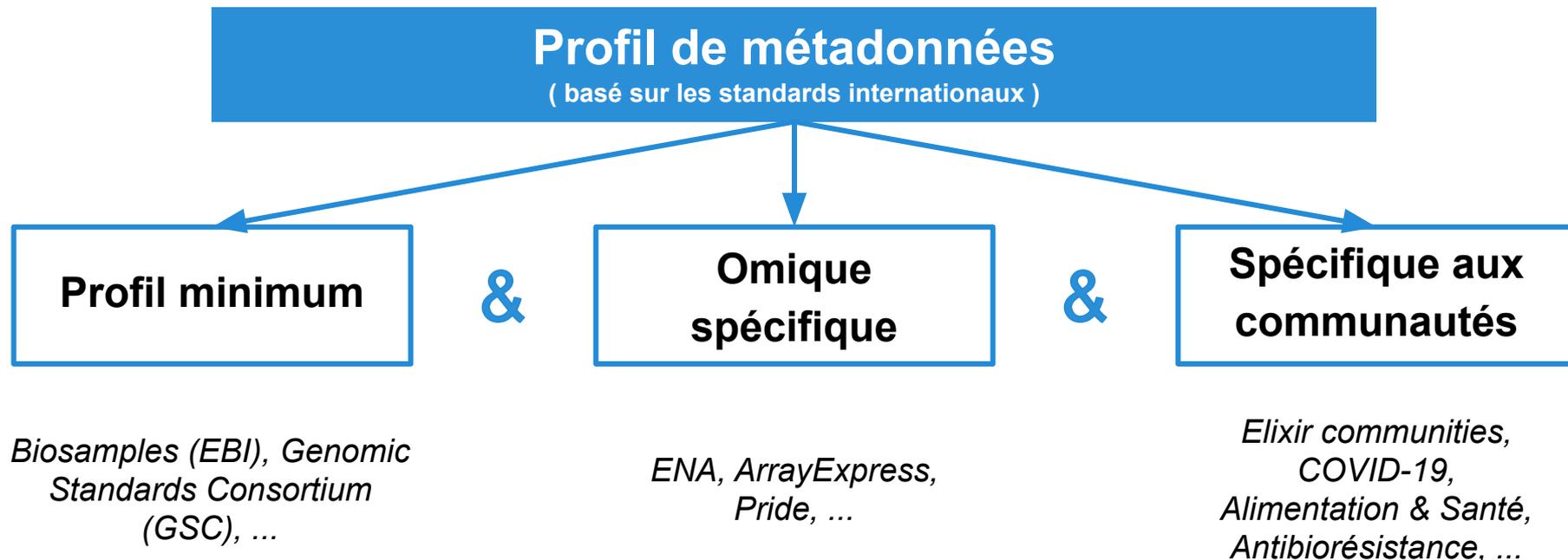
1. Identifier et collecter les informations obligatoires pour décrire et soumettre les données
2. Centraliser ces informations pour les soumettre dans une base de données adéquate
3. Implémenter des outils automatisés permettant la fluidification du flux de données

### Qui ?



(Cochrane, 2018)

## 1. Identifier et collecter les informations obligatoires pour décrire et soumettre les données



## Les métadonnées

### Définition

“Metadata are the in-depth, **controlled description** of the sample that your sequence was taken from. Essentially, the ‘**what, where, how, and when**’ of your study from collection to sequence generation, plus contextual data such as environmental conditions (latitude, longitude, temperature) or clinical observations.”

EMBL-EBI



...

### Des bases de données de métadonnées

**Biosamples**

EMBL-EBI 



**ENA**  
European Nucleotide Archive



**Array express**

**Biosample**



**GEO**  
Gene Expression Omnibus



**Sequence Read Archive**

## 2- Centraliser ces informations pour les soumettre dans une base de données adéquate

Données

Métadonnées



**ifb**  
INSTITUT FRANÇAIS  
DE BIOINFORMATIQUE

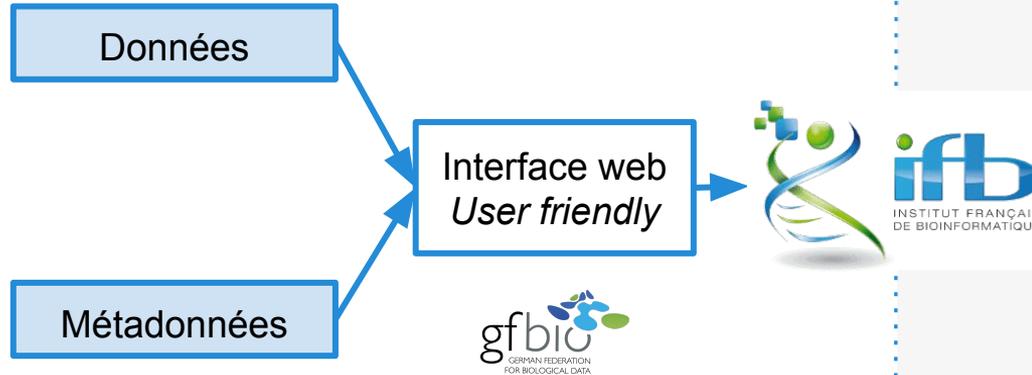
## 3- Implémenter des outils automatisés permettant la fluidification du flux de données

 **ENA**  
European Nucleotide Archive

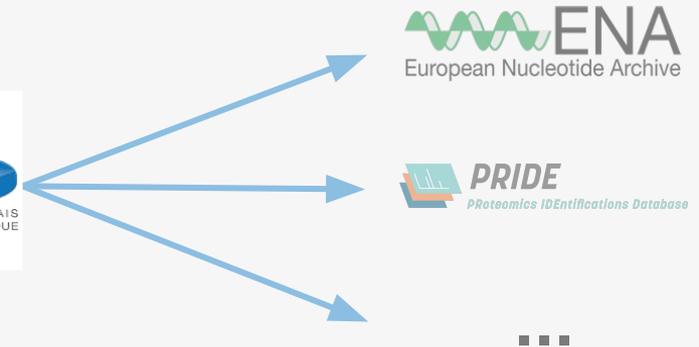
 **PRIDE**  
Proteomics IDentifications Database

...

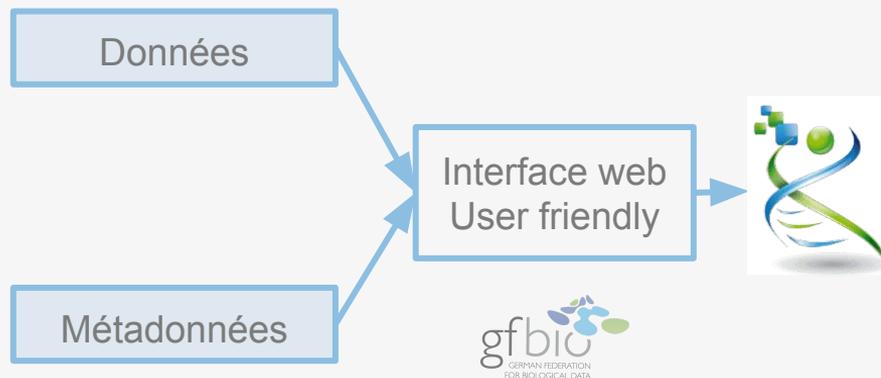
## 2- Centraliser ces informations pour les soumettre dans une base de données adéquate



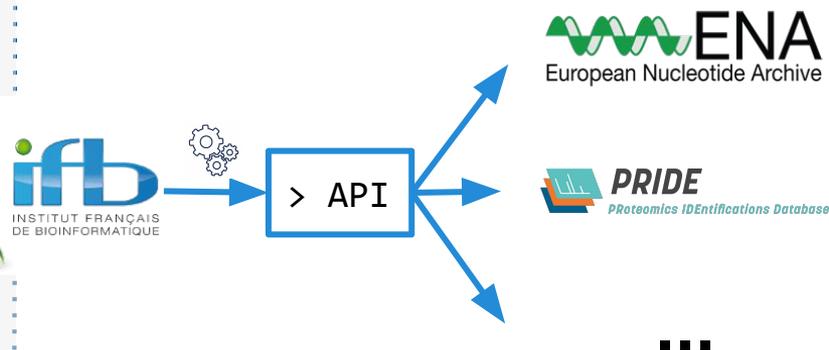
## 3- Implémenter des outils automatisés permettant la fluidification du flux de données



## 2- Centraliser ces informations pour les soumettre dans une base de données adéquate



## 3- Implémenter des outils automatisés permettant la fluidification du flux de données



## *Exemples de brokers*

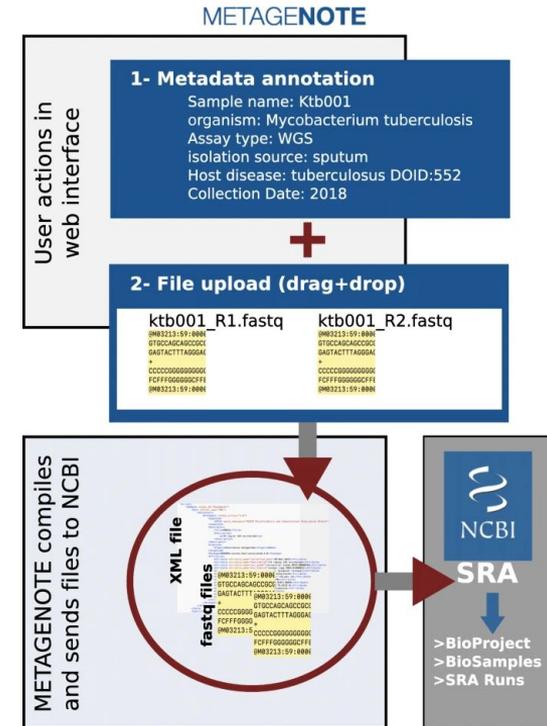
## METAGENOTE

Basé sur les *guidelines* du Genomic Standards Consortium (GSC) et des ontologies (ENVO, FMA, ...)

Interface web + MariaDb

Soumission automatique sur SRA via API

(Quiñones et al. - 2020 - <https://doi.org/10.1186/s12859-020-03694-0>)





*Preuve de concept*

## COVID-19

### Constat

SCIENCES - MEDECINE

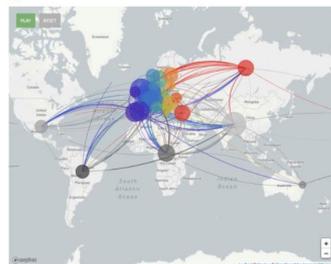
#### Covid-19 : les chercheurs français peu partageurs des séquences génétiques

La mise en commun massive permet une étude plus précise du virus et de son évolution, mais les scientifiques français y sont réticents.

Par David Larousserie - Publié le 31 août 2020 à 08h30 - Mis à jour le 01 septembre 2020 à 12h47

Lecture 7 min.

Article réservé aux abonnés



Représentation des origines des diverses importations du coronavirus en Europe entre le 7 avril et le 1er juillet 2020, tirée du séquençage de leurs génomes. Image extraite du site Noxtrain

Il n'y a de pire aveugle que celui qui ne veut pas voir. En matière de Covid-19, le dicton s'appliquerait-il à la France ? Notre pays semble en effet peu enclin à utiliser un outil de pointe qui permettrait de répondre à des questions importantes sur l'épidémie, comme

Partage



### De nouveaux outils pour faciliter la soumission

#### New tool simplifies the submission of SARS-CoV-2 data to open databases

ELIXIR Belgium and ELIXIR Germany (de.NBI) help researchers share FAIR COVID-19 data

ELIXIR Belgium, in collaboration with ELIXIR Germany and the European COVID-19 Data Platform, have developed a tool to simplify the submission of viral sequencing data to the European Nucleotide Archive (ENA), an ELIXIR Core Data Resource providing open access to nucleotide sequences. The new submission tool offers an easy-to-use interface, guides researchers through the submission process and verifies the data format and description.

Why submit data to ENA?



(Larousserie, Le Monde, 31 août 2020)

(ELIXIR, 17 November 2020)

## *Conclusion & Perspectives*

## Les actions à venir

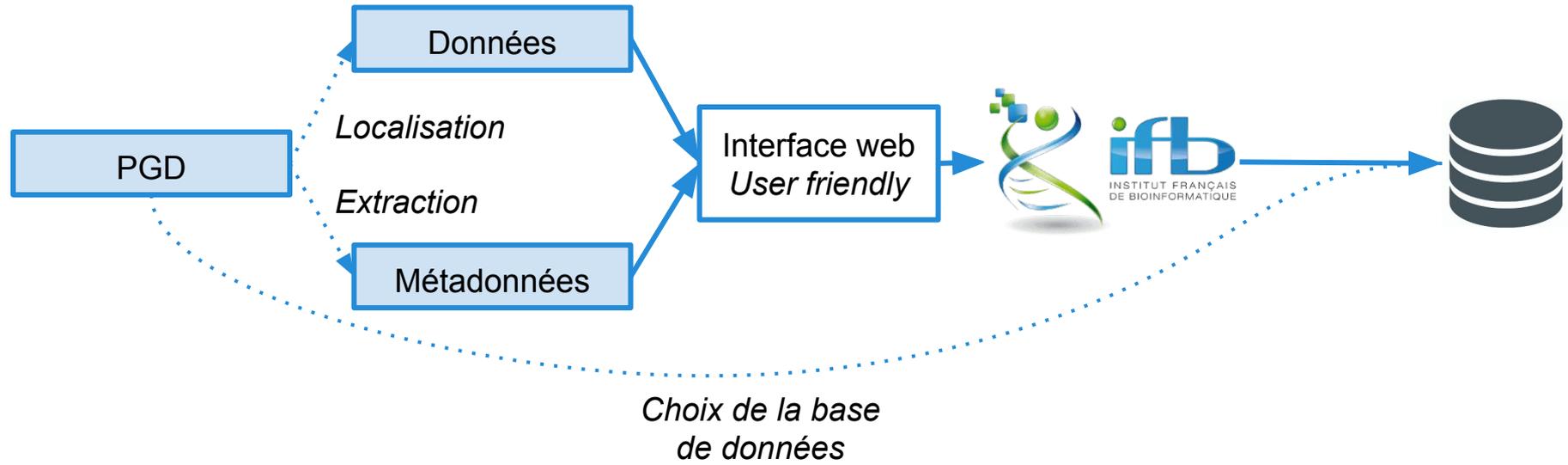
### Entre l'utilisateur et le *data broker*

- 1) Établir et collecter un ensemble d'informations nécessaires lors d'une soumission et proposer un **profil de métadonnées**
- 2) Mettre en place un **système *User friendly* de collecte** des métadonnées et des données

### Entre le *data broker* et la base de données

- 3) Mettre en place des **workflows automatiques** de soumission adaptés au type de données

## Une connexion avec le plan de gestion de données



**Mise à jour de la soumission lors de la mise à jour du PGD  
(intégration continue)**



# *Annexes*

## *Un point sur les métadonnées*

## Quelques définitions

“Metadata refers to **descriptive information** about the overall study, individual samples, all protocols, and references to processed and raw data file names.”

GEO

“Metadata characterize biological resources by **core information** including a name, a description of its input and its output (parameters or format), its address, and various additional properties.”

Encyclopedia of Database Systems, 2009

“Metadata are the in-depth, **controlled description** of the sample that your sequence was taken from. Essentially, the ‘what, where, how, and when’ of your study from collection to sequence generation, plus contextual data such as environmental conditions (latitude, longitude, temperature) or clinical observations.”

EMBL-EBI

## Genomic Standards Consortium (GSC)



**Checklist**  
(obligatoire)

&

**Environmental packages**  
(informations complémentaires)

Fichier excel téléchargeable sur GitHub  
(dernière MAJ 26 Aug 2019 )

The screenshot shows the GitHub repository page for 'GenomicsStandardsConsortium / mixs-legacy'. The repository is on the 'master' branch with 1 branch and 1 tag. It has 64 stars and 19 commits. The file list includes:

- github/ISSUE\_TEMPLATE
- 2009
- 2010
- 2011
- mixs5
- pre-2009
- gitignore
- LICENSE
- MixS\_v4.xls
- MixS\_v4.xlsx
- README.md

The README.md file is selected and shows the following content:

**mixs-legacy**

Older versions of spreadsheet formatted MixS standard templates (up to version 5).

Note that before version 4, MixS was not released as a single standard, but rather as a set of checklists. Those lists are available here in the folders labeled pre-2009, 2009, 2010, 2011.

Older versions of MixS and checklists were released as .xls files. Those files are available here, along with a new version saved in the open source .xlsx format. Some of the oldest files are available as .docx and .xslx (in the pre-2009 folder).

To request changes to the MixS standards, please use the [issue tracker](#) in the main MixS repository.

## Genomic Standards Consortium (GSC)



### **Checklist** - Informations minimales pour les séquences (MIxS)

94 items ( $\pm$  obligatoires) décrits dans 22 colonnes :

*The MIxS checklist*

*Column 1 - structured comment name: name of a checklist item as it will appear in GenBank structured comments*

*Column 2 - item: full name of item as it appears in the publication*

*Column 3 - definition: a description of the item, including links to ontologies and other resources that can be used to fill in values for the item*

*Column 4 - expected value: short description and/or expected value of an item*

*Column 5 - value syntax: the proper syntax for writing the value for a given item*

*Column 6 - example: examples of values for an item*

*Column 7 - section: the section of an item*

*Columns 8 through 18-migs\_eu,migs\_ba,migs\_pl,migs\_vi,migs\_org,me,mimarks\_s,mimarks\_c,misag,mimag,miuvig: information about whether an item is mandatory (M), conditional mandatory (C), optional (X), environment-dependent (E) or not applicable (-) for a given checklist type*

*Column 19 - preferred units: a unit suggestion if a measurement value is given*

*Column 20 - occurrence: indicates whether a given item can be used only once (1), multiple times (m), or none (0)*

*Column 21 - position: position of item as it appears in the publication*

*Column 22 - MIXS ID: a unique of an item*



## Genomic Standards Consortium (GSC)



### Environmental packages - 17 différents

Environment	Structured comment name	Package item	Definition	Expected value	Value syntax	Example	Required	Preferred unit	Occurrence	Position	MXS ID
air	alt	altitude	Altitude is a term used to identify heights of objects such as airplanes, space shuttles, rockets, atmospheric balloons and heights of places such as atmospheric layers and clouds. It is used to measure the height of an object which is above the earth's surface. In the context, the altitude measurement is the vertical distance between the earth's surface above sea level and the sampled position in the air.	measurement value	{float} [unit]	100 meter	M	meter	1		MXS:00009
air	elev	elevation	Elevation of the sampling site is its height above a fixed reference point, most commonly the mean sea level. Elevation is rarely used when referring to points on the earth's surface, while altitude is used for points above the surface, such as an aircraft in flight or a spacecraft in orbit.	measurement value	{float} [unit]	100 meter	C	meter	1		MXS:00009
air	barometric_press	barometric pressure	Force per unit area exerted against a surface by the weight of air above that surface.	measurement value	{float} [unit]	5 millibar	X	millibar	1		MXS:00009
air	carb_dioxide	carbon dioxide	Carbon dioxide (gas) amount or concentration at the time of sampling.	measurement value	{float} [unit]	410 parts per million	X	micro mole per liter, parts per million	1		MXS:00009
air	carb_monoxide	carbon monoxide	Carbon monoxide (gas) amount or concentration at the time of sampling.	measurement value	{float} [unit]	0.1 parts per million	X	micro mole per liter, parts per million	1		MXS:00009
air	chem_administration	chemical administration	List of chemical compounds administered to the host or site where sampling occurred, and when (e.g. antibodies, a herbicide, or flour); can include nucleic compounds. For chemical entities of biological interest ontology (chdbi) (v 103), <a href="http://purl.bioontology.org/bioontology/chdbi">http://purl.bioontology.org/bioontology/chdbi</a>	CHDBI:timestamp	{term:[id]} {term:[id]} {term:[id]}	ager:JC4EBE:2009:2016-05-11T20:02	X	m	1		MXS:00075
air	humidity	humidity	Amount of water vapor in the air, at the time of sampling.	measurement value	{float} [unit]	25 gram per cubic meter	X	gram per cubic meter	1		MXS:00010
air	methane	methane	Methane (gas) amount or concentration at the time of sampling.	measurement value	{float} [unit]	1900 parts per billion	X	micro mole per liter, parts per billion, parts per million	1		MXS:00010
air	misc_param	miscellaneous parameter	Any other measurement performed or parameter collected, that is not listed here.	parameter name, measurement value	{text} {float} [unit]	Biomass concentration concentration:2076 micro mole per kilogram	X	m	1		MXS:00075
air	organism_count	organism count	Total cell count of any organism (or group of organisms) per gram, volume or area of sample, should include name of organism followed by count. This method was used for the enumeration (e.g. qPCR, flow, etc.) (Should not be provided). (Example: 1000 cells per mL, qPCR)	organism name, measurement value	{text} {float} [unit] {SPCR:ATP:PF00100}	total planktonics:5.0e7 cells per milliliter:qPCR	X	number of cells per cubic meter, number of cells per milliliter, number of cells per cubic centimeter	1		MXS:00010
air	oxygen	oxygen	Oxygen (gas) amount or concentration at the time of sampling.	measurement value	{float} [unit]	600 parts per million	X	milligram per liter, parts per million	1		MXS:00010
air	oxy_sat_satp	oxygenation status of sample	Oxygenation status of sample	enumeration	{enumeration} {enumeration} {other}	aerobic	X	m	1		MXS:00075
air	perturbation	perturbation	Type of perturbation, e.g. chemical administration, physical disturbance, etc., associated with perturbation regimen including how many times the perturbation was repeated, how long each perturbation lasted, and the start and end time of the entire perturbation period. Can include multiple perturbation type.	perturbation type name, perturbation interval and duration	{text} {float} {start} {interval} {end} {duration}	antibiotic: addition:R12018-05-11T14:30Z:2016-05-11T19:30Z:P1430M	X	m	1		MXS:00075
air	pollutants	pollutants	Sublist types and amount or concentrations measured at the time of sampling; can report multiple pollutants by entering numeric values preceded by name of pollutant.	pollutant name, measurement value	{text} {float} [unit]	lead:0.15 microgram per cubic meter	X	gram, mole per liter, milligram per liter, microgram per cubic meter	1		MXS:00010
air	resp_part_matter	respirable particulate matter	Concentration of substances that remain suspended in the air, and comprise mixtures of organic and inorganic substances (PM10 and PM2.5); can report multiple PMs by entering numeric values preceded by name of PM.	particulate matter name, measurement value	{text} {float} [unit]	PM2.5:10 microgram per cubic meter	X	microgram per cubic meter	1		MXS:00010
air	samp_safety	sample safety	Safety is the total concentration of all dissolved salts in a liquid or solid (in the form of an extract obtained by centrifugation) sample. While safety can be measured by a complex chemical analysis, this method is official and less concerning. More often, it is instead derived from the conductivity measurement. This is known as practical safety. These derivations compare the specific conductance of the sample to a safety standard such as seawater.	measurement value	{float} [unit]	1 milligram per liter	X	milligram per liter, practical safety unit, percentage	1		MXS:00010
air	samp_store_dur	sample storage duration	Duration for which the sample was stored.	duration	{duration}	P19M	X	m	1		MXS:00011
air	samp_store_loc	sample storage location	Location at which sample was stored, usually name of a specific free-air room.	location name	{text}	Freezer no.5	X	m	1		MXS:00075
air	samp_store_temp	sample storage temperature	Temperature at which sample was stored, e.g. -80 degrees Celsius.	measurement value	{float} [unit]	-80 degree Celsius	X	degree Celsius	1		MXS:00011
air	samp_vol_wt_or_ext	sample volume or weight for DNA extraction	Volume (mL), weight (g) of processed sample, or surface area swabbed from sample for DNA extraction	measurement value	{float} [unit]	1500 milliliter	X	milliliter, gram, milligram, square centimeter	1		MXS:00011
air	solar_irradiance	solar irradiance	The amount of solar energy that arrives at a specific area of a surface during a specific time interval.	measurement value	{float} [unit]	1.36 kilowatts per square meter per day	X	kilowatts per square meter per day, ergs per square centimeter per second	1		MXS:00011
air	temp	temperature	Temperature of the sample at the time of sampling.	measurement value	{float} [unit]	25 degree Celsius	X	degree Celsius	1		MXS:00011
air	ventilation_rate	ventilation rate	Ventilation rate of the system in the sampled premises.	measurement value	{float} [unit]	750 cubic meter per minute	X	cubic meter per minute, liters per second	1		MXS:00011
air	ventilation_type	ventilation type	Ventilation system used in the sampled premises.	ventilation type name	{text}	Operable windows	X	m	1		MXS:00075
air	volatile_organic_compounds	volatile organic compounds	Concentration of carbon-based chemicals that easily evaporate at room temperature; can report multiple volatile organic compounds by entering numeric values preceded by name of compound.	volatile organic compound name, measurement value	{text} {float} [unit]	formaldehyde:500 nanogram per liter	X	microgram per cubic meter, parts per million, nanogram per liter	1		MXS:00011
air	wind_direction	wind direction	Wind direction is the direction from which a wind originates.	wind direction name	{text}	Northwest	X	m	1		MXS:00075
air	wind_speed	wind speed	Speed of wind measured at the time of sampling.	measurement value	{float} [unit]	21 kilometer per hour	X	meter per second, kilometer per hour	1		MXS:00011



BioSamples stocke et fournit des descriptions et des métadonnées sur les échantillons biologiques utilisés dans la recherche et le développement par les universités et l'industrie.



## Sample

Sample content reference:

Field	Description	Type	Cardinality =====
name	The short name of the sample.	String	Required
release	The date at which the sample was first made public.	Date ISO 8601	Required
update	The date at which the sample was last updated.	Date ISO 8601	System Generated
domain	The AAP domain the sample belongs to.	String	Required
accession	The sample unique identifier in the BioSamples database. If not provided, one will be automatically assigned.	String	Required for [PUT] requests
characteristics	The key-value pairs representing the attributes of the sample.	Object	Optional
externalReferences	A list of links towards external references, such as datasets in other archives.	Array	Optional
relationships	A list of relationships this sample has to other, existing, samples.	Array	Optional
data	A more structured data format to allow submission of tables (eg. antibiogram) in addition to key-value pairs.	Array	Optional, required only for structured data submission using our POST, PUT or PATCH endpoints.  Note - You must provide an AAP domain for your data. It can be same as the sample domain if you are the submitter of both the sample metadata and structured data.

[https://www.ebi.ac.uk/biosamples/docs/references/api/submit#\\_submission\\_minimal\\_fields](https://www.ebi.ac.uk/biosamples/docs/references/api/submit#_submission_minimal_fields)





## La qualité actuelle des métadonnées

### **The variable quality of metadata about biological samples used in biomedical experiments**

Rafael S. Gonçalves & Mark A. Musen

Scientific Data volume 6, Article number: 190021 (2019)

11.4 millions de métadonnées testés issues de BioSample (NCBI) et BioSamples (EBI)

### **Bilan**

- Les noms et les valeurs ne sont pas contrôlés ni standardisés
- Les valeurs ne sont pas toujours du bon type (ex : binaire)
- Manque d'outils de validation qui éviteraient des aberrances