

Curation and Quality Control of Data in EpiFlu™ Database

Anne Pohlmann
FLI Riems

Singapore
2nd BLE-GISAID Symposium
24/25 October 2014



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

The host of GISAID database

As of 2011 the Federal Republic of Germany represented by the Federal Ministry of Agriculture and Food (BMEL) is official host of the GISAID database

- **Friedrich-Loeffler-Institute (FLI), Federal Research Institute for Animal Health**
 - Content-based data curation
 - Data import and correction
 - Contact point for questions concerning data curation
- **Federal Office for Agriculture and Food (BLE)**
 - Technical operation and maintenance of the database and web server
 - Monitoring and surveillance of worldwide availability
 - Data security and integrity



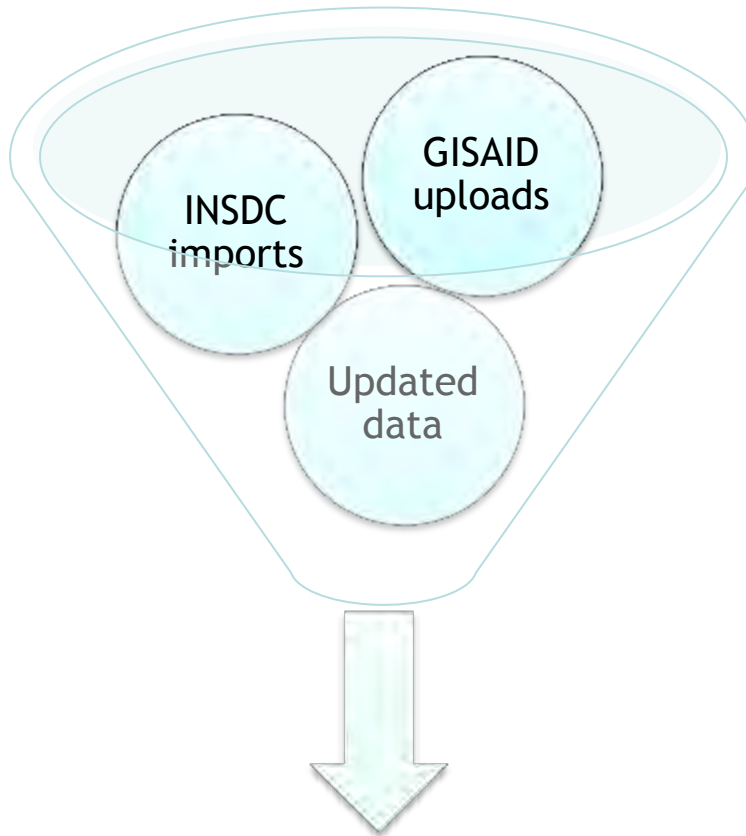
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

What data need to be curated?



Curated Data

Nucleotide sequences

- Length, code
- Coding sequences
- Protein sequences
- ...

Virus information

- Type, subtype, lineage
- Dates
- Host information
- Geographic information
- Resistance information
- Antigenic information
- ...



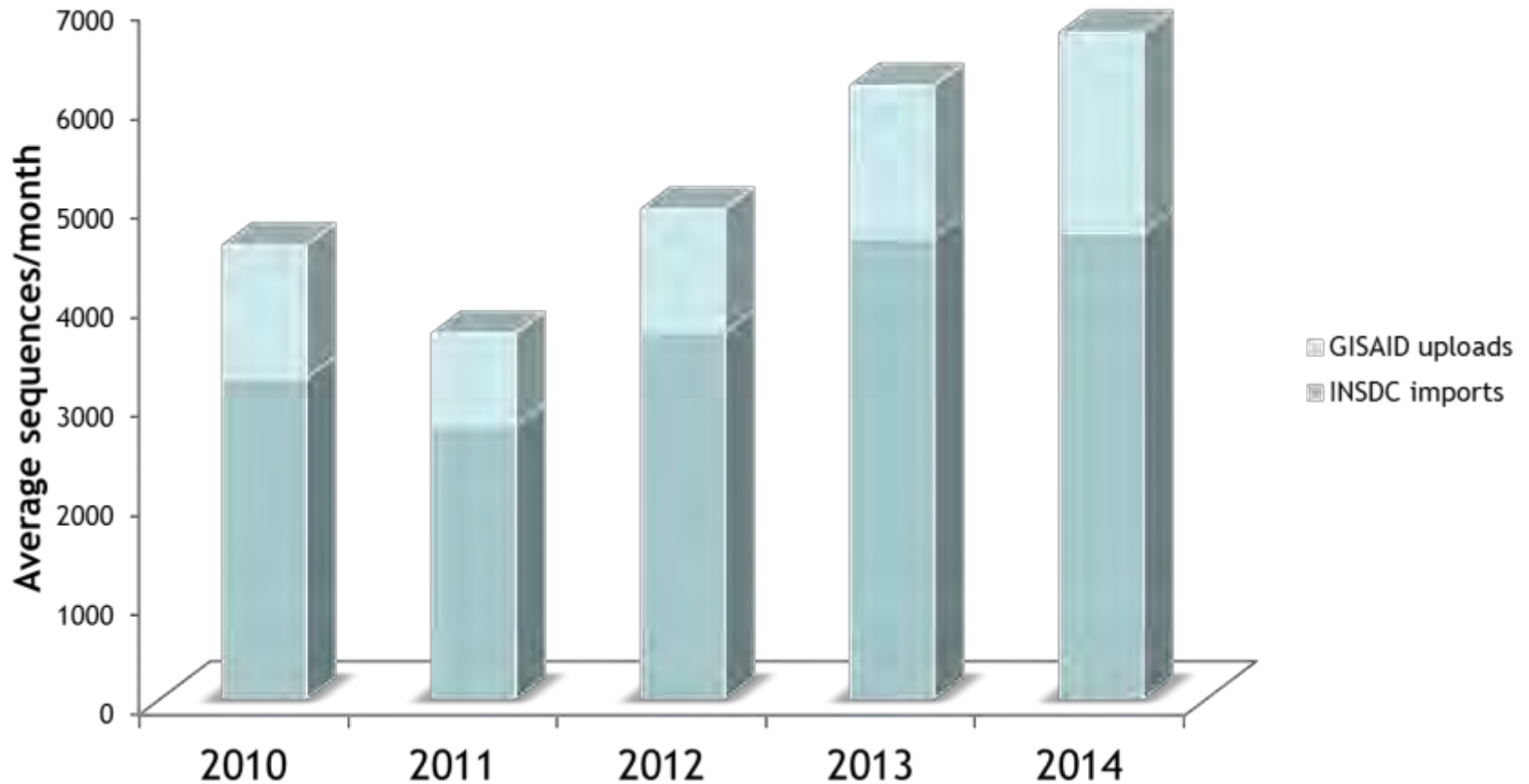
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

How many data needs to be curated?



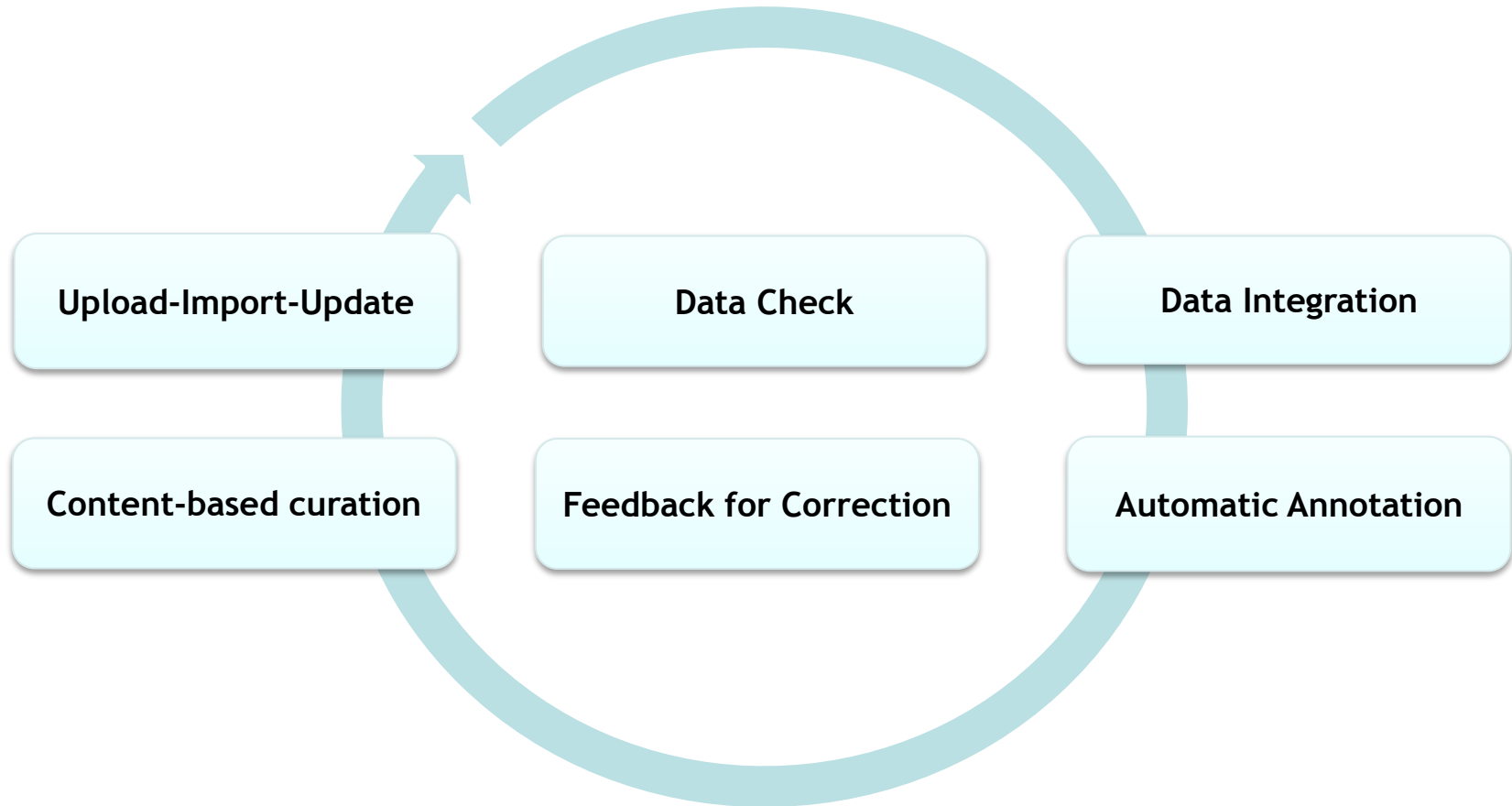
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Data curation is an interactive and iterative process



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Upload tools assists submitter to ensure correct and complete data

Data upload:

- Single Upload
- Batch upload



Submitter

GISAID GISAID published: 4.269 viruses with 8.130 Sequences Total count: 104

► Browse ► **Upload** ► Workset Management ► Administration

Virus

Virus Detail

Virus Name *
Example: A/Wisconsin/2145/2001 or A/chicken/Rostov/864/2007

Passage Category

Passage History *
Example: c1/c2

Harvest Date

Virus Type * Subtype H Subtype N Lineage

Host Group *
None
Animal
Environment
Human
Laboratory derived
Unknown

None
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17



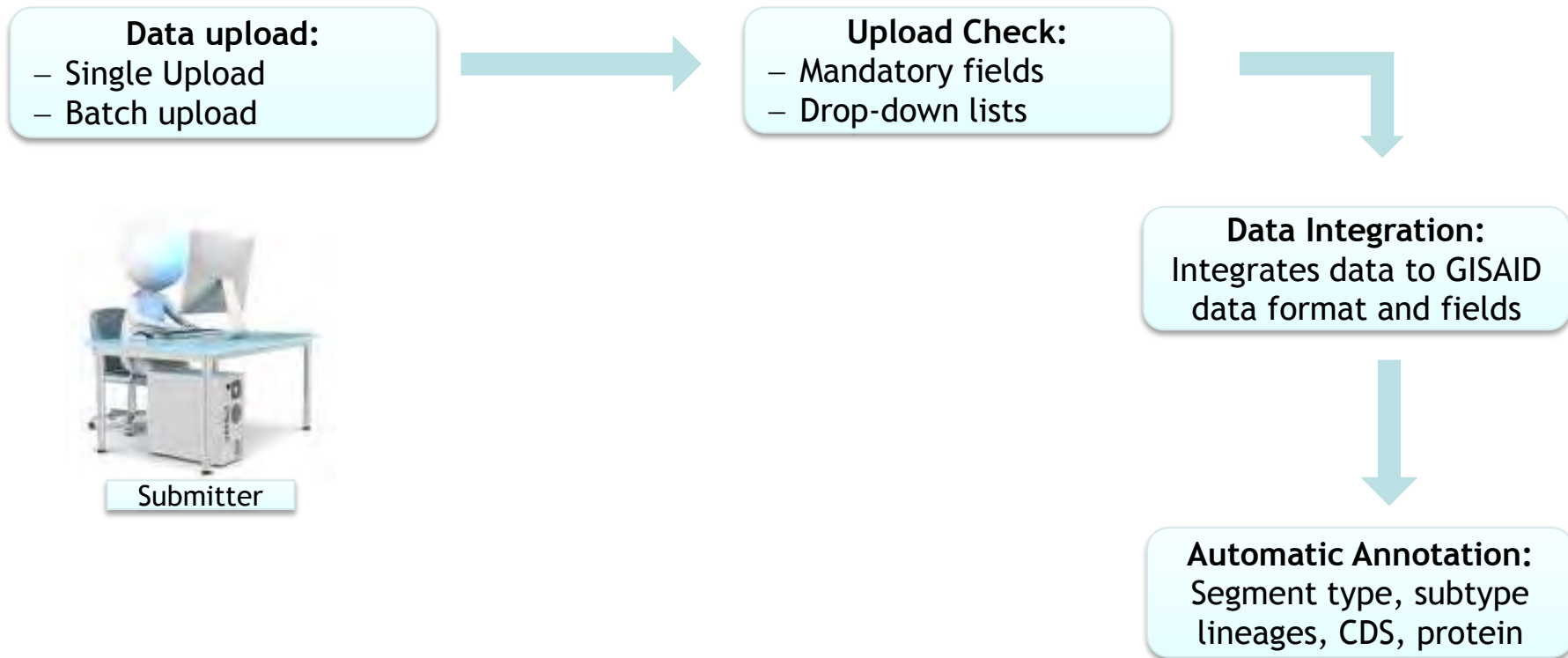
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Interactive process to improve data quality



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Automatic annotation in GISAID 2.0



CLC Qiagen annotation server:

- Automatic annotation of segment type, influenza type and lineages
- Identification of virus subtype from HA and NA segments
- Determination of coding sequences and integration of corresponding proteins

CDC Atlanta H5N1 Clade annotation:

- Designation of clades for H5N1 sequences



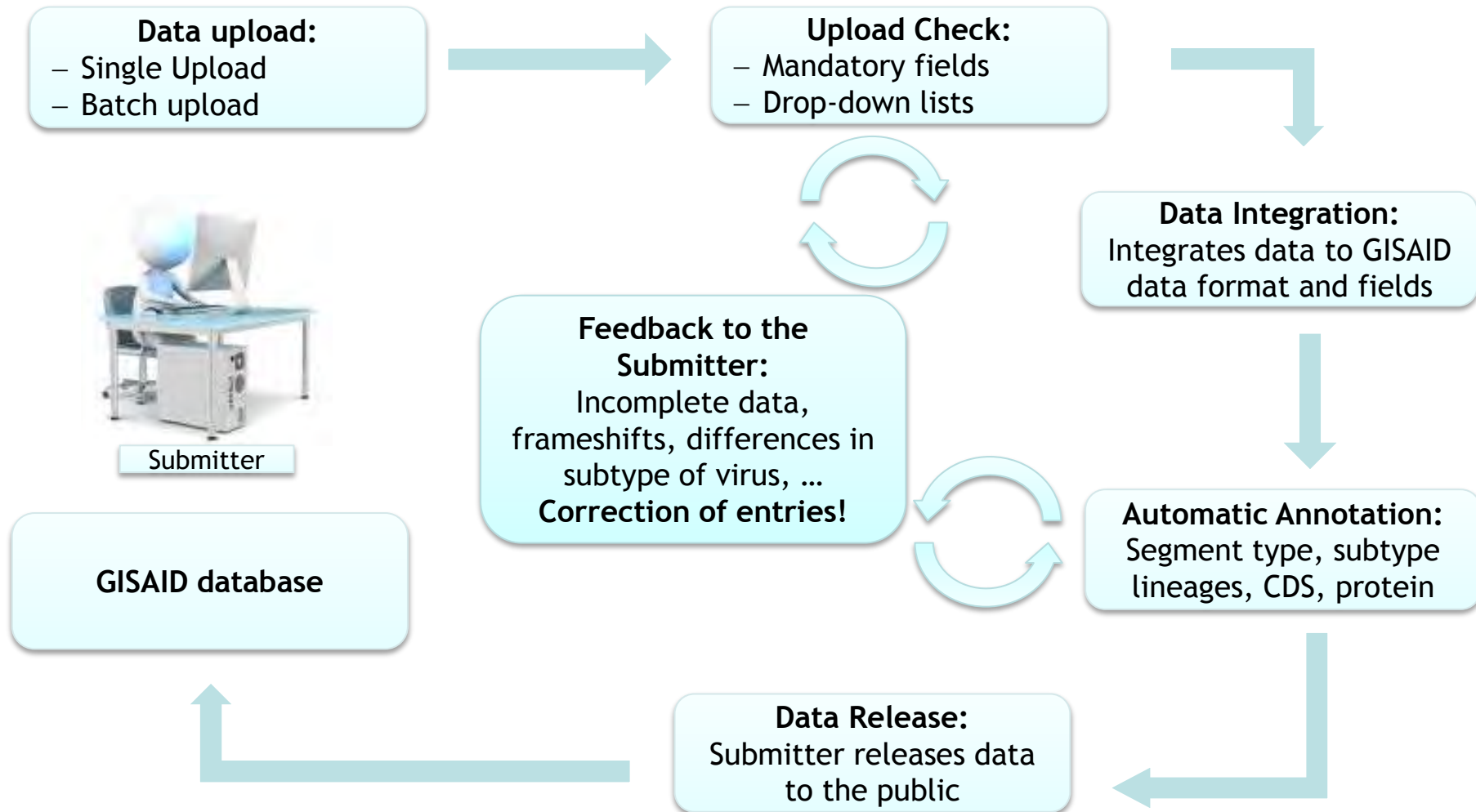
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Interactive process to improve data quality of GISAID uploads



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Import and quality control of INSDC data

Data import:

- New Influenza data
- Updated influenza data



Import Check:

- Parse data
- Check data



Curator

Correction by Curator:
Complete data,
identification of passages,
check for duplication,
insert detailed location,
use scientific host name
(Avian), annotation, ...

GISAID
GISAID published: 4,269 viruses with 8,130 Sequences Total count: 104,611 viruses with 75,300 Sequences

Import EBI-Data

Change selected virus Change host Change country Delete selected Delete all

Search for:
☐ Influenza A ☐ Influenza B ☐ Influenza C ☐ last month

EBI-Query:

☐ Show all updated ☐ Show all new

Accession	Updated	New	Virus Type ID	Tax ID	Name	Passage	Collect date	Country / State	City	Submission
-----------	---------	-----	---------------	--------	------	---------	--------------	-----------------	------	------------



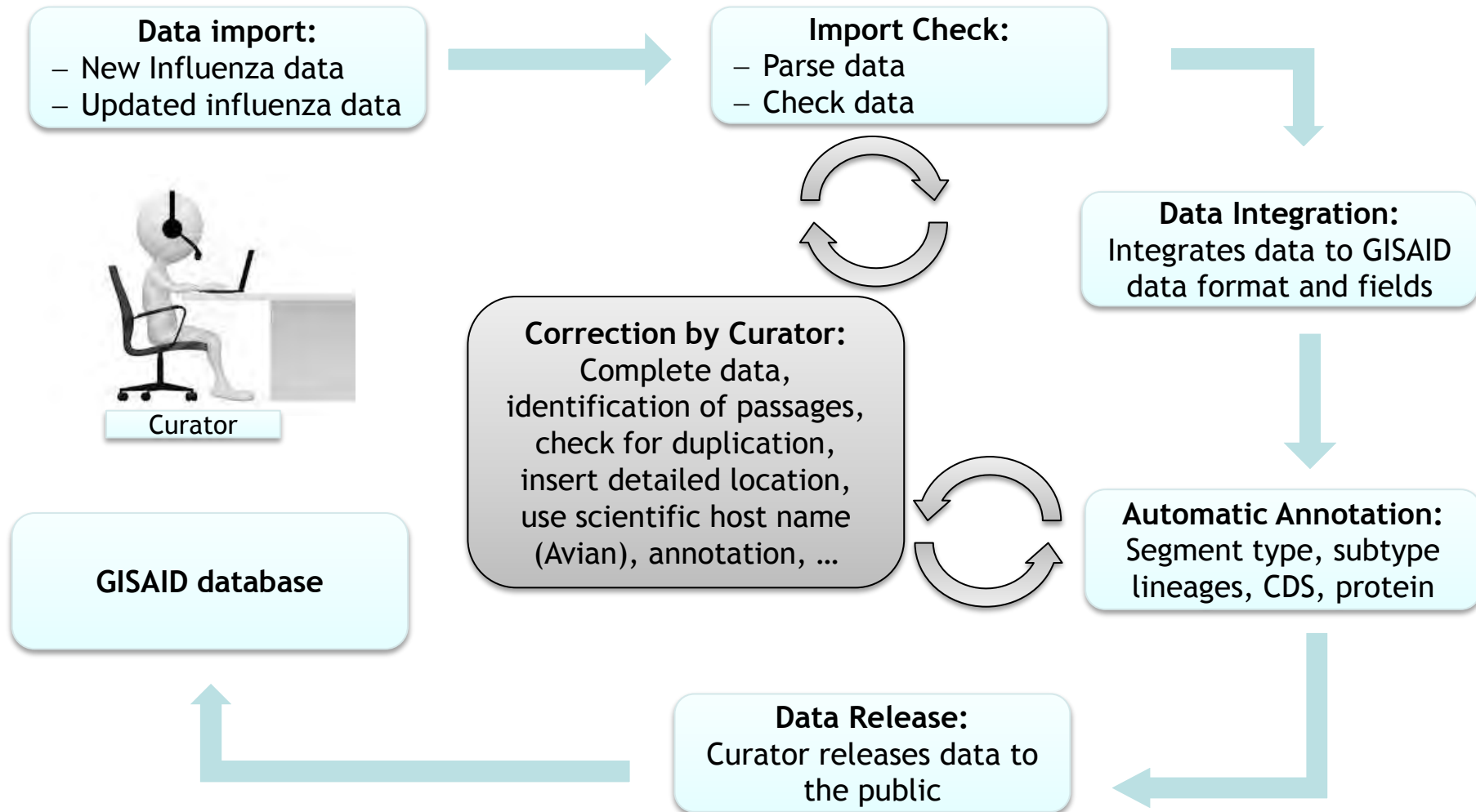
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Import and quality control of INSDC data



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Curation of Catalogue Content



GISAID database

GISAID

GISAID published: 4,269 viruses with 8,130 Sequences Total count: 104,611 viruses with 75,300 Sequences

[Browse](#) [Upload](#) [Workset Management](#) [Administration](#) [Registered Users](#)

Host Species

Host Group ▲	Host Subgroup ▲	Host Family ▲
Animal	Avian	Chicken
Environment	mammals	Duck
Human		Eagle
Laboratory derived		Falcon
Unknown		Goose
		Grouse
		Guineafowl
		Gull
		Ostrich
		Other avian
		Partridge
		Passerine
		Pheasant
		Sandpiper
		Shearwater
		Swan
		Turkey
		US Quail

Administration

- My Profile
- Users
- Catalogs**
 - Statistics
 - Import Viruses
 - System Configuration
- Animal Domestic Statuses
- Animal Health Statuses
- Animal Specimen Sources
- Countries
- Drugs
- Geographic Groupings
- Host Families
- Host Groups
- Host Species
- Host Subgroups
- Host Subspecies
- Human Outbreaks
- Human Patient Statuses
- Human Specimen Sources
- Virus Types
- Institutions
- Lineages
- Passage Categories
- Resistance Levels
- Rights
- Roles
- Segment Types
- Protein Types



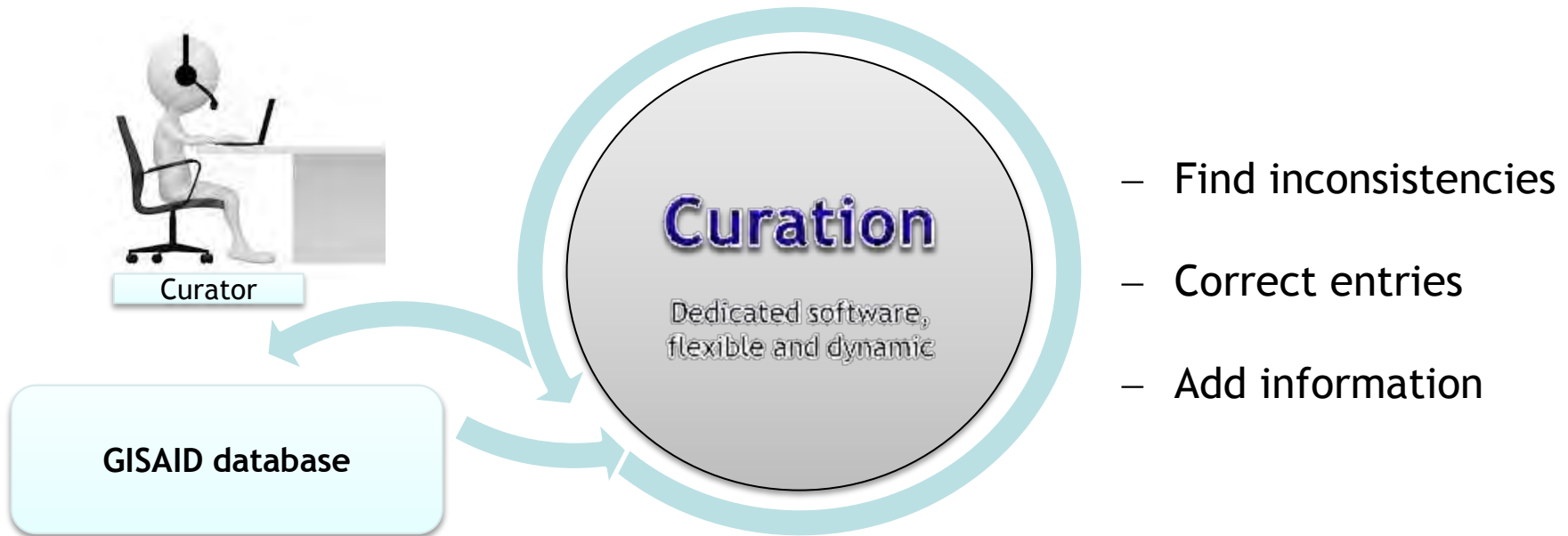
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Content-based data curation



Isolate_Id	Isolate_Name	Subtype	Lineage	Location	Host	Collection_Date	Update_Date
EPI_ISL_90291	A/Turkey/Germany/R617/2007	A / H6N2		Europe / Germany	Turkey	2011-05-26	2011-05-26
EPI_ISL_90292	A/Turkey/Ontario/6118/1968	A / H8N4		North America / Canada	Turkey	2011-05-26	2011-05-26
EPI_ISL_94864	B/MACAU/200392/2011				Human	2011	
EPI_ISL_94865	B/MACAU/601485/2010	B	Victoria		Human	2010-08-26	2011-08-15

Lack of location

Inconsistent collection date



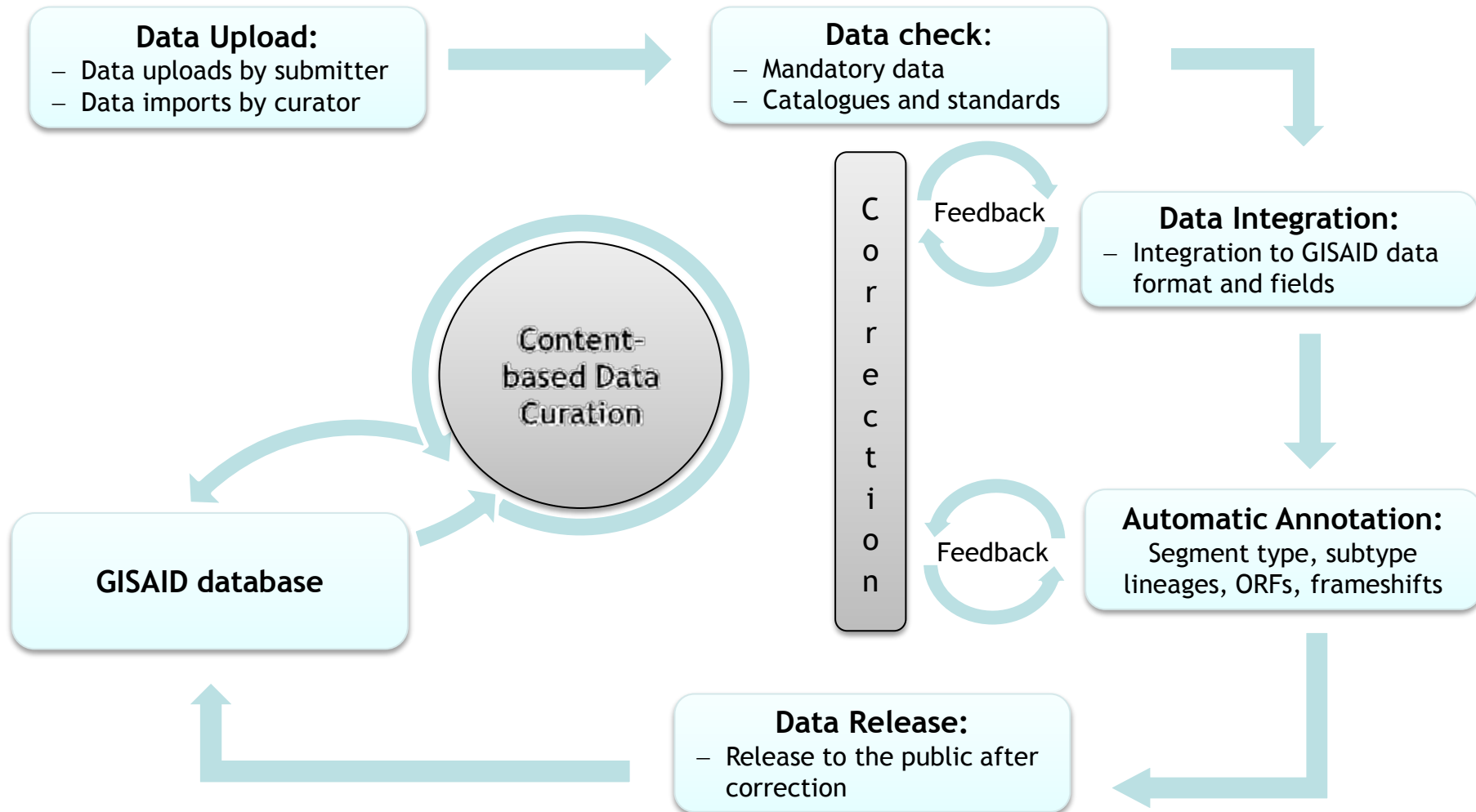
FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

The GISAID 2.0 Curation Concept



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health

Thanks

FLI | Martin Beer | Dirk Hoeper | Thomas Mettenleiter

BLE | Jan Mark Pohlmann | Holger Heuser | Hanns-Christoph Eiden

GISAID | Alan Hay | Peter Bogner

CDC Atlanta | Catherine Smith | Rebecca Garten | Nancy Cox

WHO CC Melbourne | Naomi Komadina

Scopeland | Claudia Peißert | Robert Lembcke | Marko Bacetic

CLC Qiagen | Ben Turner

BII | Sebastian Maurer-Stroh | Raphael Lee Tze-Chuen



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health



Thank you for your attention.



service@gisaid.org

curator@gisaid.org



FRIEDRICH-LOEFFLER-INSTITUT

since 1910

FLI

Bundesforschungsinstitut für Tiergesundheit
Federal Research Institute for Animal Health