



INFRASTRUCTURE FOR

LONG-TERM PRESERVATION AND

OCR ANALYSIS OF HERBARIUM IMAGES

ICEDIG.EU

Innovation and consolidation for large scale digitisation of natural heritage



Full text search portal





≥ Ш

Ш Т

S

< Σ

F

()

Ω

S

R

ш

> 4 million specimens have been ingested at CINES computing center based in Montpellier, France

OPEN

The open data portal CINES enables a full text search on metadata and OCR results to display a list of matching results and then to access the herbarium sheet. URL : https://opendata.cines.fr

File format

Checksum

JPEG/1.01

SHA-256

543184a5fd534e5e9155b3bc5a54273b8d0c3fe8f06c1ef69a9576a3a3299f

and excels (minist) off experiences, Ad Proto COMMON BERLOF Parent of Justice do Bioardian on 12-57 2 COT1 or analyze 27-5000 Dia 2011 nnhn > Hubert Basis of record StillImage|PRESERVED_SPECIMEN Institution code MNHN http://coldb.mnhn.fr/catalognumber/mnhn ocurrence Id p00184040 Family Asteraceae Hubertia humblotii (Klatt) C.Jeffrey cientific name Labat, J.N. Yahaya, I. Daroueche, E. Djoubeiri, M. Recorded by Event date 1999-11-19 Field notes Arbuste 2 m assez fréquent à ces altitudes: capitule petit à 3-4 fleurs ligulées; ligule jaune, petite; fleur tubulaire absente ou peu nombreuse (photo) Location Comoros grande como Convalescent MNHN mind adency chiving date 2018-09-4 Archiving Id ark:/87895/1.90-21503 P00184040.jp File name

ownload OCR resi

d'altitude Fourre dense has 2.3 m presque WII/m monospecifique a Phillipi Lithosol sur basalte Arbuste 2 m assez frequent a ces altitudes capitule petit a a 4 fleurs ligulees ligule jaune petite fleur tubulaire absente ou peu nombreuse (photo Phenologie fl adoublesa B CNDRS G K MO 19/11/1999 Labat J N 3167 aVec Yaha)a I , Daroueche E & Djoubeiri M

ш in herbaria all over the world. Most scientific corpora is CHALL available in analog form but physical copies are fragile. Digital copy must be preserved without loss of information in future.

- global trend to industrial digitizing
- data difficult to handle even for medium size institutes same challenges being faced by potentially hundreds of herbaria in Europe
- makes sense to work together to develop a solut

impacts on research activity

- full text search on OCR results
- sharing knowledge on image analysis, specimens will be
- discoverable by the entire scientific community
- characterisation of features, annotation, meta data extraction
- organisation of the data and metadata in certain formats data curation to ensure that meta data and data formats remain meaningful in future



the whole processing chain including, quality control, OCR analysis on HPC facilities and indexing on a full text searching using Elastic Search engine has been implemented.

metadata

- Darwin Core and Dublin Core based
- \cdot have been imported in the Trusted Digital Repository (TDR) and mapping has been done on the long-term archival system from the Global Biodiversity Information Facility



GBIF Global Biodiversity Information Facility



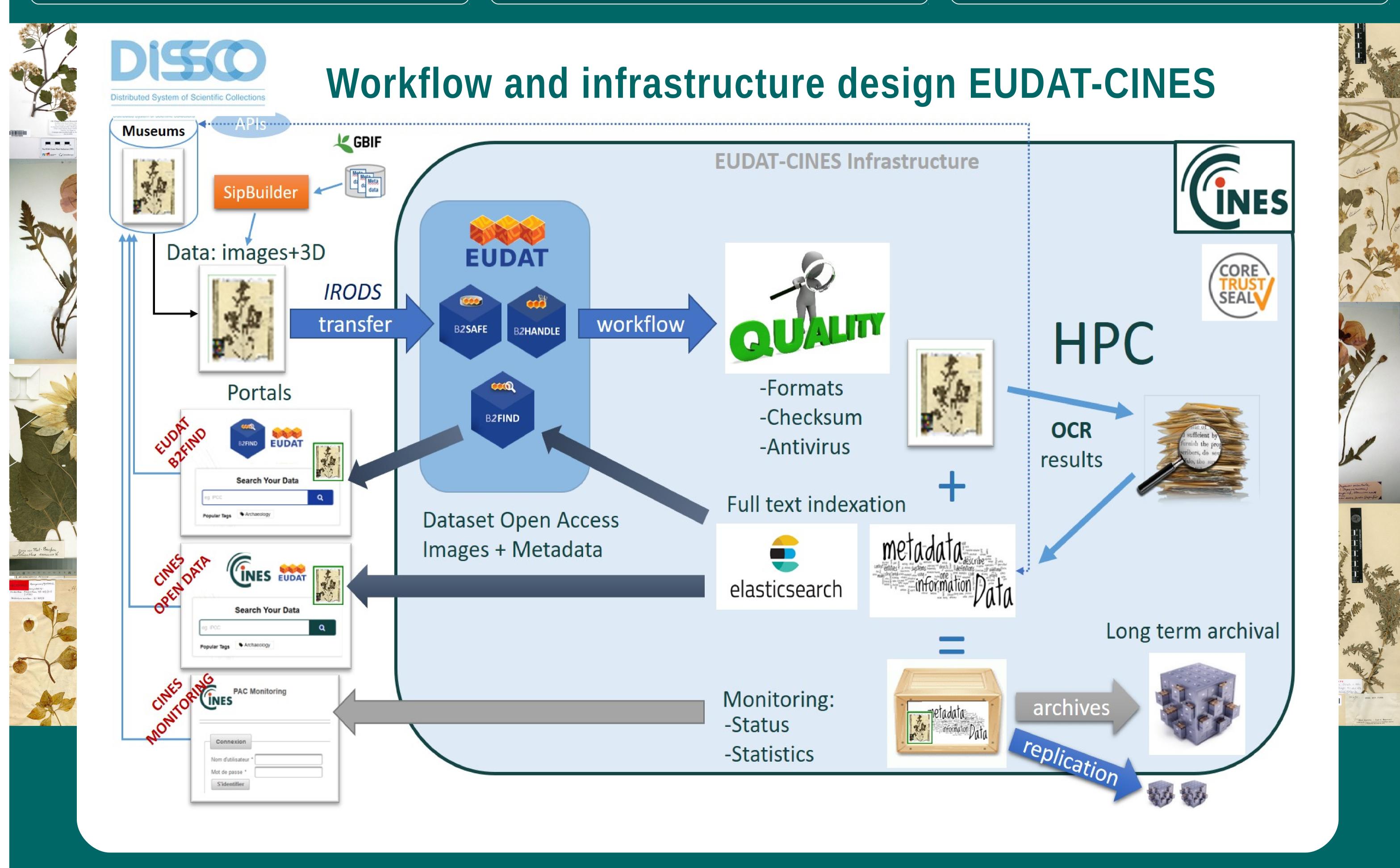
eTDR European Trustworthy Digital Repository

- \cdot securing and certifying long-term archiving
- \cdot long distance replication with EUDAT partners.



scientific use-cases

- mining for duplicates in collections located in different museums
- \cdot exploring the data using deep learning algorithms
- identifying seldom species



ICEDIG – "Innovation and consolidation for large scale digitisation of natural heritage" - is an EU-funded project that aims at supporting the implementation phase of the new Research Infrastructure DiSSCo ("Distributed System of Scientific Collections") by designing and addressing the technical, financial, policy and governance aspects necessary to operate such a large distributed initiative for natural sciences collections across Europe.