# Arkindex, an open-source platform for document image recognition applied to archeology

Christopher Kermorvant, TEKLIA, Paris

Mélodie Boillet, TEKLIA, Paris

Pablo Ciezar, INRAP

Anne-Violaine Szabados, CNRS - ArScAn
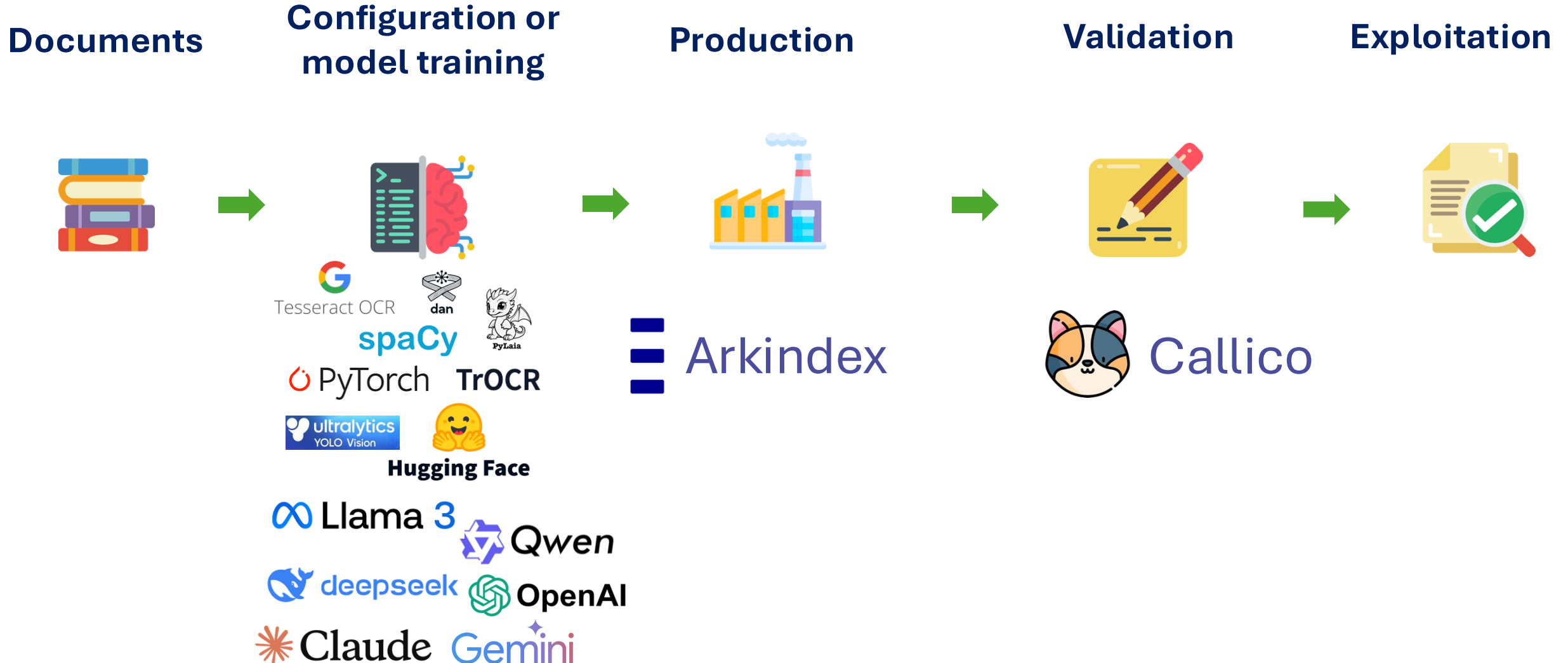
Julien Schuh, UPN - Université Paris Nanterre

# Automatic document processing with AI

AI models can recognize, understand and extract data from historical documents

The challenges are now:

- How to choose the right model: accurate, easy to run and cheap

- How to model data coming from different sources (meta-data, human annotation, predictions from AI)

- How to process documents at large scale: reliability, completeness, quality control
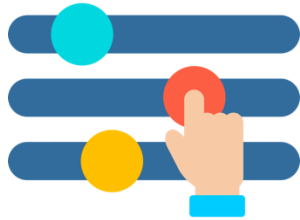
# Document processing workflow

**Documents**     **Configuration or model training**     **Production**     **Validation**     **Exploitation**

# Arkindex: a generic document processing platform

Developed by TEKLIA since 2019
Already used for more than 60 projects

## Customization

Process all types of documents

## Scaling up

Process 1000 or 10 million pages

## Open-source

Dissemination and community participation

https://gitlab.teklia.com/arkindex

# Arkindex: a generic data model

# Arkindex: a generic data model

# Arkindex: a generic data model



Manual and automatic classification

# Arkindex: a generic data model



Manual and automatic text transcription
(printed and handwritten)

# Arkindex: a generic data model



Meta-data integration

# B. Bruyère's handwritten excavation journals

Journals of excavations by Egyptologist Bernard Bruyère from 1922 to 1951 at Deir el-Medina, on the left bank of Thebes

A lot of information, drawings of objects and surveys of monuments were not included in the published reports and are only available in the handwritten journals.

# B. Bruyère's handwritten excavation journals

Automatic transcription (HTR) of 1000 handwritten pages

Manual illustration detection and annotation (metadata)

Indexing and search

# Identification of Argonne roller-stamped motifs

Late Antique Argonne ware was produced on a massive scale in Northwestern Europe between the 4th and 6th centuries.
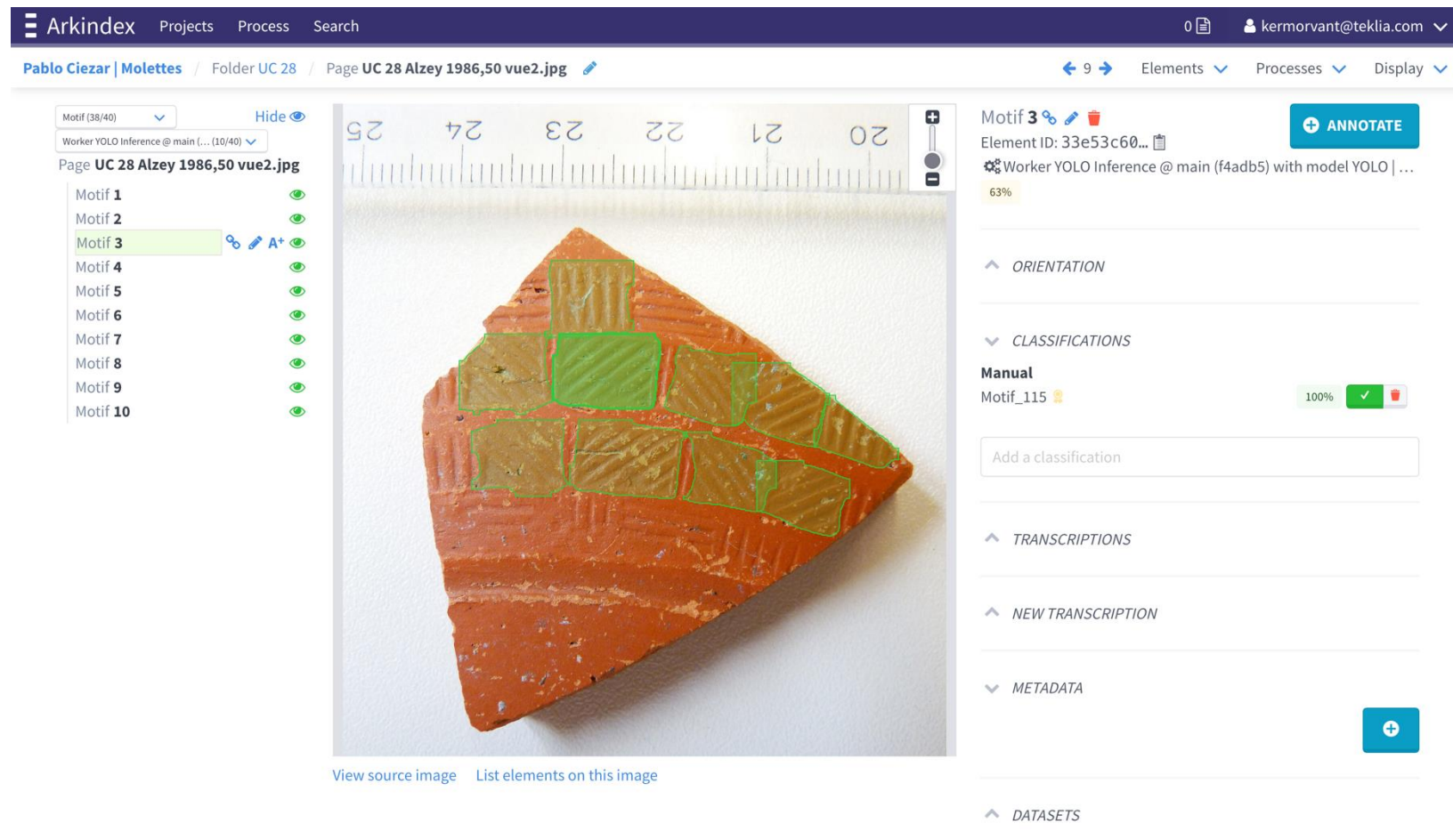
The roller-embossed decoration allows for the identification of production sites and workshops, and provides elements of dating.

# Identification of Argonne roller-stamped motifs

<u>Detection of the motifs</u> : image segmentation algorithm (YOLO), stored as image element

<u>Type of motif</u> : classification on the detected element

# Information extraction from Celtic coin catalogues

A vast amount of information on coins is available in paper publications.

These publications contain descriptive, geographical and bibliographical information in semi-structured form.

# Information extraction from Celtic coin catalogues

[OCR and Structured information extraction](#) using multimodal Large Language Models (vLLM)

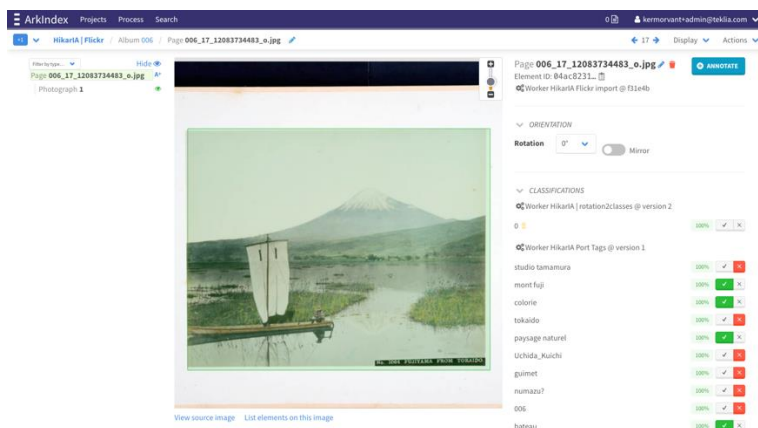# Scaling-up processing



Processing workflow definition



Distributed processing
(multiple servers, HPC)

# TEKLIA's open-source software suite for document processing
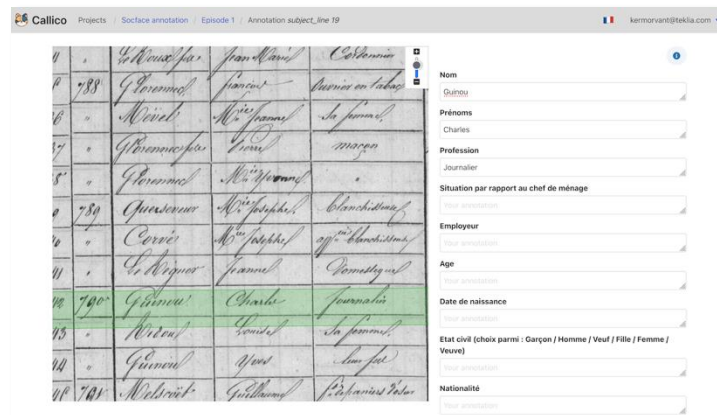
## Arkindex
Document processing



Open-source and Enterprise licences

https://support.teklia.com

## Callico
Collaborative annotation campaigns



Open-source

https://gitlab.teklia.com

## HuggingFace
Models and datasets



Open-source, open-weights

https://huggingface.co/Teklia

kermorvant@teklia.com