
TEXT-BASED INFORMATION EXTRACTION SERVICES

Dr. Eduard Santamaria

eduard.santamaria@iosb.fraunhofer.de

Berlin, 1-2 June 2017



Fraunhofer

IOSB

TEXT-BASED INFORMATION EXTRACTION SERVICES

- Given a piece of text (not an image) with the contents of the specimen label, extract important metadata about the specimen:
 - Scientific name
 - Collector
 - Collection date
 - Geo-coordinates
 - Location

MUSCI AUSTRALASIAE EXSICCATI
Edited by H. Streimann
FISSIDENTACEAE
373. *Fissidens linearis* Brid. var. *linearis*, Musc. Rec. Suppl.
4: 187. 1819.
AUSTRALIA. Western Australia: Boyagin Rock, Boyagin Nature Reserve, 20 km NW of Pingelly. 32° 28'S 116° 53'E, alt. 350 m.
Large exposed prominent rock outcrop surrounded dry sclerophyll forest.
On ground on boulder shaded by a larger boulder.
12 September 1994 H. Streimann 54200
Mus. Bot. Berol.

i

Annotation Types

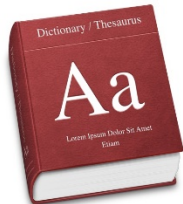
<input checked="" type="checkbox"/> Date	<input type="checkbox"/> Documen...	<input checked="" type="checkbox"/> GeoCoords	<input checked="" type="checkbox"/> Location	<input checked="" type="checkbox"/> Person
<input checked="" type="checkbox"/> Scientific...				

MAIN PROBLEMS

- Errors in text (OCR does not always work)
 - Deterioration of specimen sheet
 - Multiple languages (with particular characters)
 - Handwritten parts
- Syntax and abbreviations
 - Authors try to fit much information in a small label
- Ambiguity
 - Stone – family name or description of the place?
 - White – family name, color?

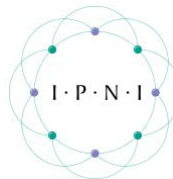
COLLECTOR

```
MUSCI AUSTRALASIAE EXSICCATI
Edited by H. Streimann
FISSIDENTACEAE
373. Fissidens linearis Brid. var. linearis, Musc. Rec. Suppl.
4: 187. 1819.
AUSTRALIA. Western Australia: Boyagin Rock, Boyagin Nature
Reserve, 20 km NW of Pingelly. 32° 28'S 116° 53'E, alt. 350 m.
Large exposed prominent rock outcrop surrounded dry sclerophyll forest.
On ground on boulder shaded by a larger boulder.
12 September 1994 H. Streimann 54200
Mus. Bot. Berol.
```



Best matches

Harvard University
Herbaria



The International
Plant Names Index

Best matches

Take into account multiple forms:
Jan-Peter Frahm
J.-P. Frahm
Frahm, J.-P.
Accompanied by Leg., Det. or
equivalent form?

RESULT

DATE AND GEO-COORDINATES

- Date and Geo-coordinates services are based on regular expressions.
 - "(\\d{1,2})\\.\\s(\\w+)\\s(\\d{4})" → 5. VII 1875
 - "(\\w+)[\\,|\\.]?\\s(\\d{1,2})[\\.|\\,]\\s+(\\d{4})" → Mayo, 21. 2000
 - "(\\-?\\d{1,3})\\s?[°d*]?\\s?(\\d{1,3})(?:\\'|\\")\\s?(\\d{1,3})(?:\\|\\'|\\")\\s?([nsewNSEW])"
→ 47° 16'39" E
 - etc.
- Multiple languages for months:
 - English, English short form, French, Spanish, German, Italian, Portuguese
- Numbers: 6, 06, vi, but also v1 and o6

USAGE OPTIONS

■ Web User Interface


Annotator Dienste

Request Form API Documentation

Herbarium label text

Annotator

Scientific Names

Find 

■ Java Library



■ REST Web Service

StanDAP-Herb Webservices ^{0.1.0}

default



POST /scientific-names

POST /dates

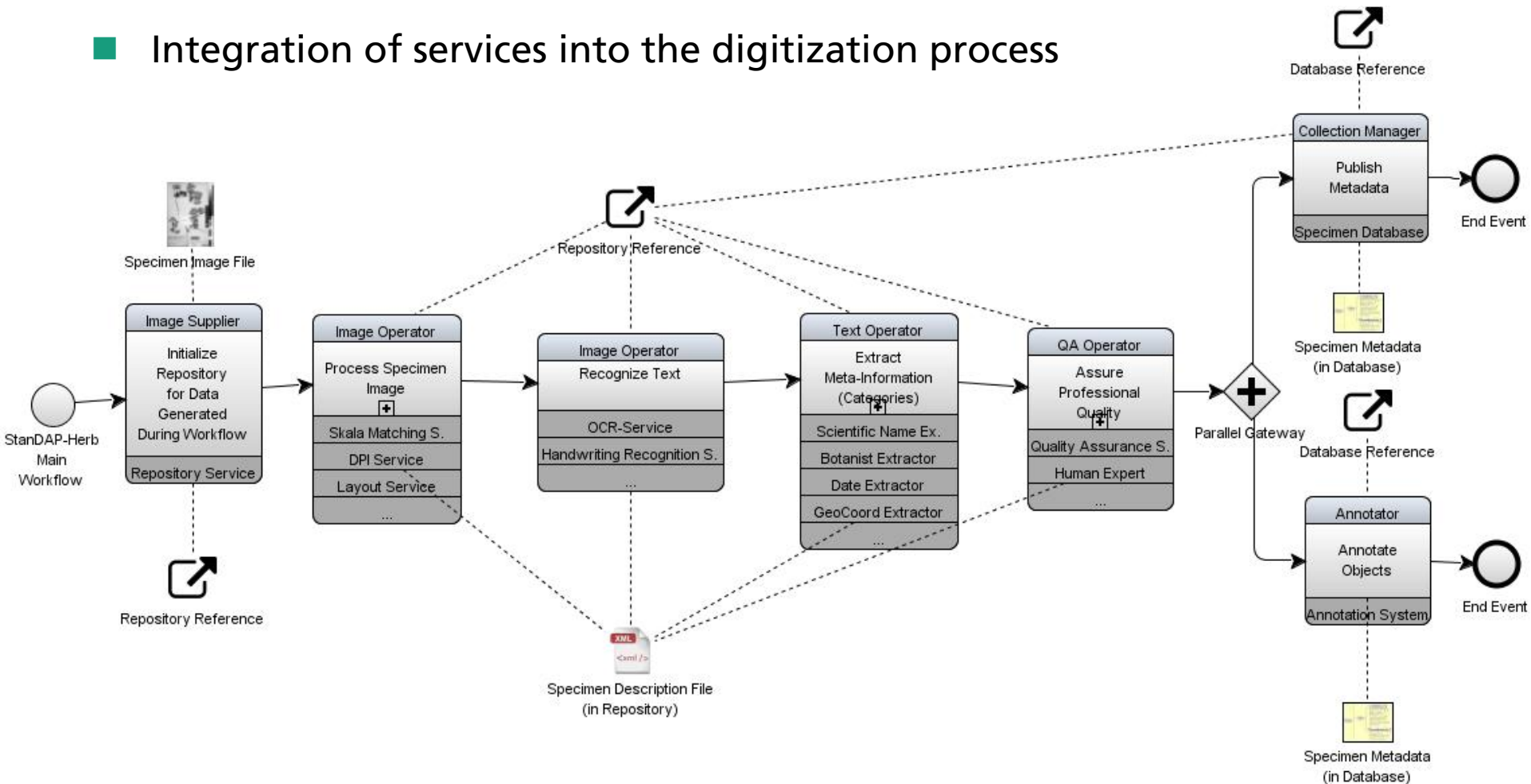
POST /botanists

POST /geo-coordinates

POST /locations

WORKFLOW

Integration of services into the digitization process



RESULTS WITH TEST DATASET (I)

1162 text files with output from OCR.

Precision & Recall values automatically computed comparing service results with known correct results.

	Precision	Recall
Scientific Name	0,50	0,48
Collector Name	0,73	0,23
Date	0,66	0,50
GeoCoordinates	0,48	0,33
Country Name*	0,85	0,54

Precision = $TP / (TP + FP)$

Recall = $TP / (TP + FN)$

* Service provides other locations
but only country was compared.

RESULTS WITH TEST DATASET (II)

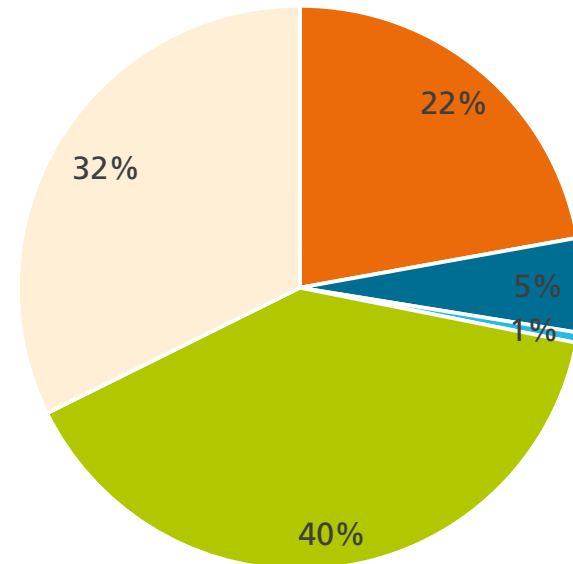
- 100 random files
- 100 plant names

If OCR'd text is ok, scientific name is found most of the time.

A few sheets contain lots of plant names



Results



- 1. Bad text quality (37)
- 2. FN: Not found in local dictionary (9)
- 3. FP1: Incorrect plant name (1)
- 4. FP2: Name on sheet, but it's not the specimen (66)
- 5. TP (54)

RESULTS WITH TEST DATASET (III)

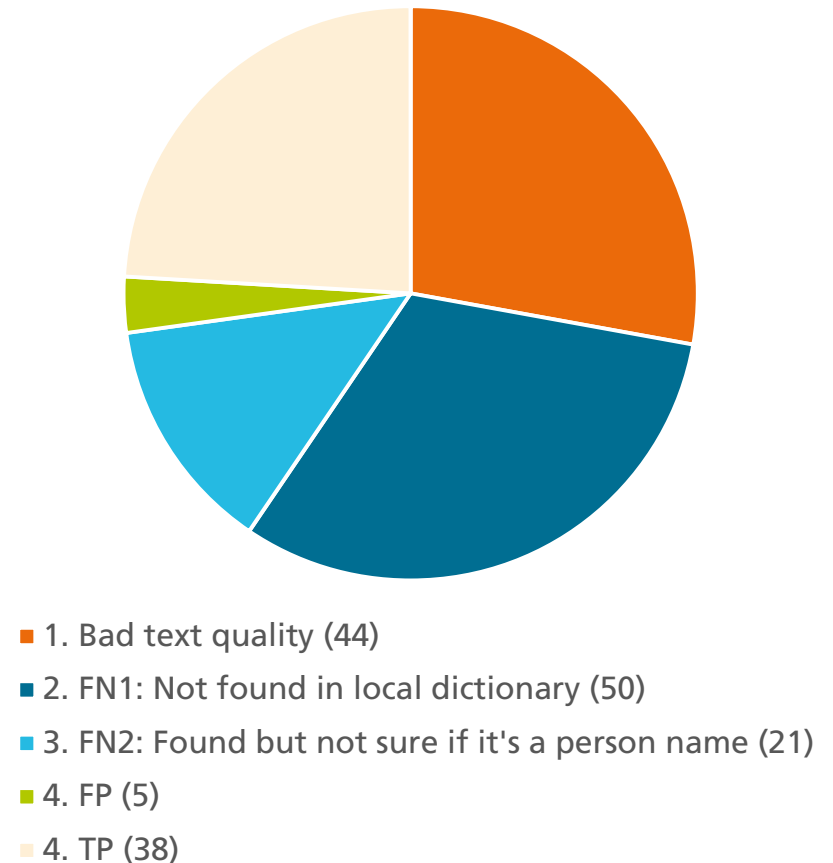
- 100 random files
- 153 collector names

Number of good results limited by:

- Text not good enough for the service.
- Extent of dictionaries *
- Limitations of the annotator or ambiguity (ex. Stone)

* Remote dictionary queried only if local hit found.

Results



POTENTIAL IMPROVEMENTS

- Automatically update local dictionaries periodically.
- Extend local dictionaries with input from users as specimen sheets are processed.
- Use statistical indicators to select candidates:
 - Which one of a number of candidates did the user more often select as the correct answer?
- Currently, the information extraction services work in isolation:
 - Date plus collector name could be checked with an external database to determine location.

