

# READ

RECOGNITION & ENRICHMENT  
OF ARCHIVAL DOCUMENTS

---

## D5.13

### Page Image Explorer (PIE) P3

---

Markus Diem, Florian Kleber  
CVL

Distribution: <http://read.transkribus.eu/>

**READ**  
**H2020 Project 674943**

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 674943



<b>Project ref no.</b>	H2020 674943
<b>Project acronym</b>	READ
<b>Project full title</b>	Recognition and Enrichment of Archival Documents
<b>Instrument</b>	H2020-EINFRA-2015-1
<b>Thematic priority</b>	EINFRA-9-2015 - e-Infrastructures for virtual re- search environments (VRE)
<b>Start date/duration</b>	01 January 2016 / 42 Months

<b>Distribution</b>	Public
<b>Contract. date of delivery</b>	31.12.2018
<b>Actual date of delivery</b>	22.11.2018
<b>Date of last update</b>	18.12.2018
<b>Deliverable number</b>	D5.13
<b>Deliverable title</b>	Page Image Explorer (PIE) P3
<b>Type</b>	report
<b>Status &amp; version</b>	in progress
<b>Contributing WP(s)</b>	WP5
<b>Responsible beneficiary</b>	CVL
<b>Other contributors</b>	CVL
<b>Internal reviewers</b>	NAF, ASV
<b>Author(s)</b>	Markus Diem, Florian Kleber
<b>EC project officer</b>	Christophe DOIN
<b>Keywords</b>	Document Clustering, Visualization

## Contents

<b>1</b>	<b>Executive Summary</b>	<b>4</b>
<b>2</b>	<b>PIE Database Tool</b>	<b>4</b>
<b>3</b>	<b>PIE</b>	<b>6</b>
<b>4</b>	<b>Outlook</b>	<b>7</b>

---

# 1 Executive Summary

The Page Image Explorer (PIE) allows intuitive exploration of documents. The key idea is to access potentially unsorted document collections and connect/group their items by user defined criteria. Hence, PIE strongly focuses on user interaction and visualization of large document collections. PIE will be built upon the READ Framework<sup>1</sup> which is publicly available under LGPLv3.

In D5.13 a tool to create the PIE database, which is the basis for the clustering, was created. This tool is integrated into the READ framework system, takes a base folder as input and generates a json file (database) with all the features as output. The prototype uses the information for the document clustering and visualizes the information. The PIE database tool and PIE itself are available as opensource on github under LGPLv3.

Section 2 describes the features used for clustering and the PIE database tool. Section 3 shows the PIE prototype and Section 4 summarize the tool and gives an outlook.

## 2 PIE Database Tool

To create a database with all features used for clustering the PIE database tool has been developed. It is part of the ReadFramework and freely available on github, see <https://github.com/TUWien/ReadFramework>. The user specifies the database path, which contains all images and the corresponding PAGE XML files. The PIE database tool crawls all PAGE XML files in the specified database path (including subdirectories) and saves all features to a json database file. The creation of the database needs about 2 minutes for about 6000 PAGE XML files.

Based on experiments features proposed in D5.12 have been skipped, while e.g. transcribed text content has been added to cluster according the text contents of document images. The following features are finally saved in the database file for each document:

- document size
- number of images in a document + image sizes
- number of tables in a document + tables sizes
- number of charts in a document + charts sizes
- number of text regions/text lines in a document + text regions/text lines sizes
- transcribed text present in the document

Additionally, a global dictionary of all appearing words together with their occurrence is stored in the json file. The Listing 2 shows an example of the json database file created by the PIE database tool:

---

<sup>1</sup><https://github.com/TUWien/ReadFramework>

```

1 {"database": "C:\\pie\\data",
2   "dictionary": {
3     "Bara": 9,
4     "Pere": 6,
5     "donsella": 4,
6     "fill": 3,
7     "habitant": 4,
8     "pages": 4,
9     "rebere": 8,
10    "texidor": 4,
11    "viudo": 3
12  },
13  "imgs": [
14    {"imgName": "0001_1551_orig_list of punishment_9.jpg",
15     "xmlName": "C:\\pie\\data\\0001_1551_orig_list of punishment_9.xml",
16     "height": 4394,
17     "width": 5475,
18     "content": " Pro al pr de d...",
19     "textRegions": [
20       {
21         "height": 1373,
22         "width": 1964
23       },
24       {
25         "height": 3316,
26         "width": 1983
27       }
28     ],
29     "imgRegions": [
30       {
31         "height": 1373,
32         "width": 1964
33       }
34     ]
35   },
36   {"imgName": "0001_1551_orig_list of punishment_10.jpg"
37     ...
38   }
39 }

```

The features have been selected based on experiments. However, more features as proposed in D5.12 could be added for certain document collections. These could be:

- paper color
- text color

- writer (based on writer identification)

The content as suggested in D5.12 has already been added to the database tool.

### 3 PIE

The PIE visualization interface, available on github <https://github.com/TUWien/PIE>, mainly shows the embedding viewport which features OpenGL for responsive zooming and panning. Figure 1 shows a picture of the UI. Here, each document is represented by a dot in the embedding space. The embedding space groups similar document pages (the user defines which similarity function(s) to take). Hence, a page's position rather relies on its relationship with all other documents than typical Euclidean distances. In Figure 1, all document pages are visualized as colored dots, where the color indicates the document group. On the right side a list with all document folders used to generate the json file shows up. The different colors indicate different documents. The left plot shows all documents according to their image width and height. The second plot shows the documents according to the average region width and height. Thus, e.g. pages with a certain region width (second plot) or a certain image size (first plot) could be clustered/selected for further exploration.



Figure 1: UI of PIE showing *all documents* with different features selected. The bottom screenshots shows the selection of a document. It can be seen that it's pages are clustered in the feature space (red dots).

---

## 4 Outlook

The PIE database tool and PIE have been developed in D5.13 to show that a clustering on selected features allow for an intuitive exploration of unsorted documents. The PIE database tool crawls an unsorted set of documents and creates a database with all extracted features. PIE can be used for visualization of the proposed features in 2D, e.g. based on the text content or e.g. to get all pages with tables. This allows a sorting and clustering of documents, thus an exploration of unsorted documents. The visualization utilize OpenGL viewports with Qt overpainting (for e.g. nice font rendering), which allows PIE to scale better compared to its precursor. Based on User Experience (UE) studies, features could be added based on specific properties of certain collections, which could be done as future work.