# CLEAT Document Image Database Overview

## 1. Executive Summary

A large, heterogeneous collection of documents is a fundamental ingredient for research on document image analysis and recognition and the development of information retrieval systems. This document provides an overview of the CLEAT document image database, produced by the Laboratory for Language and Media Processing (LAMP) at the University of Maryland, which contains:

- Multi-lingual binary and grayscale images scanned directly from University of Maryland handwriting collections

- Binary and grayscale images scanned directly from unstructured multi-lingual document image sources, including business correspondence, technical journals and handwritten notes

- Binary and grayscale images scanned directly from structured multi-lingual document image sources, including dictionaries, phonebooks and yellow pages

- Vast amount of images automatically cropped from on-line sources

- Binary images scanned from 1st and other generation photocopies of real forms

- Document images assembled from existing document databases

- All document images tagged with page level attributes for each page

- Software for viewing document images, visualize and edit ground truth information

- Document segmentation ground truth generated on select subset of document images, with each segment zoned and tagged

- Signature and logo detection ground truth manually created on select subset of document images, with each ground truth region zoned and tagged

- Signature and logo detection results generated by algorithms developed by University of Maryland on select subset of document images, with each detected region zoned and tagged

To accommodate diverse needs in document image analysis and understanding research, the CLEAT document image database includes images with a rich blend of imaging resolutions and degradations.

This volume of the database contains 16,854 images, utilizing a total of more than 6.0 gigabytes of storage. Table 1 provides a detailed description of the CLEAT database in terms of genre type and language distributions.

**Table 1: Genre types and distribution of the CLEAT document image database.**

| Forms, Drawing, Tables et at. | |
|---|---|
| Forms | 644 |
| Drawing | 42 |
| Tables | 100 |
| Chemistry formulae | 25 |
| Math equations | 165 |
| Figures | 40 |
| Total | 1016 |

| Business documents and Memo letters | |
|---|---|
| Business documents clean | 52 |
| Business documents degraded | 2700 |
| Business documents with annotations | 160 |
| Memo letters (English + Multilingual) | 978 |
| Total | 3890 |

| Journal and Conference Papers, Articles | |
|---|---|
| English | 2785 |
| German | 359 |
| Japanese | 478 |
| Total | 3622 |

| Newsletters and Flyers | |
|---|---|
| Google images | 1417 |
| Arabic Newswire + Broadcast News | 338 |
| Total | 1755 |

| Structured Documents | |
|---|---|
| Phonebook | 229 |
| Dictionaries (Chinese English, English Chinese) | 1148 |
| Yellowpages | 84 |
| Total | 1461 |

| Handwritten | |
|---|---|
| Arabic | 60 |
| Chinese | 146 |
| Cyrillic | 410 |
| Japanese | 47 |
| Korean | 80 |
| Thai | 319 |
| Hindi | 281 |
| Total | 1343 |

| Page Classification Datasets | |
|---|---|
| Document | 797 |
| Image with Text | 1695 |
| Non-Document | 1275 |
| Total | 3767 |

| | |
|---|---|
| Total in all genre categories | 16854 |

## 2. Page-level Information

Page-level attributes provide essential information associated with each document image. They are manually created for the entire CLEAT database, and can be further edited using the GEDI software provided.

**Table 2: Attribute and value sets used in the ground truth of CLEAT database.**

| Attributes | Value 1 | Value 2 | Value 3 | Value 4 | Value 5 | Value 6 | Value 7 | Value 8 |
|---|---|---|---|---|---|---|---|---|
| ImageClass | Document | ImageWithText | Non-Document | | | | | |
| PageType | Printed | Handwritten | Mixed | | | | | |
| GenreClass | Business-Memo Letter Form Other | Article-Table Figure Other | Graphics-Drawing | Structured-Phonebook Yellowpage Dictionary Other | Newsletters | Other-BroadcastNews Newswire | Unknown | |
| PrimaryLanguage | English | German | Chinese | Japanese | Korean | Thai | Cyrillic | Arabic |
| Source | Tobacco | Web | ScannedMedia | Other | | | | |
| Quality | Good | Poor | | | | | | |
| Misc | | | | | | | | |

Table 2 lists the set of common attributes and their associated values used in the creation of CLEAT page-level ground truth. Figure 2 shows an example XML ground truth file displayed using browser.

```
<?xml version="1.0" encoding="UTF-8" ?>
<!-- GEDI is developed at Language and Media Processing Laboratory, University of Maryland. -->
- <GEDI xmlns="http://lamp.cfar.umd.edu/GEDI" version="1.0">
  - <DL_DOCUMENT src="aah97e00-page02_2.tif" NrOfPages="1" docTag="xml">
      <DL_PAGE gedi_type="DL_PAGE" src="aah97e00-page02_2.tif" pageID="1" width="2592" height="3300" Quality="Good" ImageClass="Document" PageType="Printed"
        Source="Tobacco" PrimaryLanguage="English" GenreClass="Business-Memo" Misc="Blank" />
    </DL_DOCUMENT>
  </GEDI>
```

**Figure 1: Visualization of page-level attributes in XML format.**

## 3. Zone-level Information

Zone-level attributes provide essential information related to a specific region on a document page. The list of zone-level attributes can be defined in an extensible fashion by an application end user. This tight integration with the end application enables training and evaluation of various document image analysis and recognition algorithms on the CLEAT database.

Figure 2 shows a list of zone types, in which the "*DL*" prefix indicates the zone types in the ground truth data.
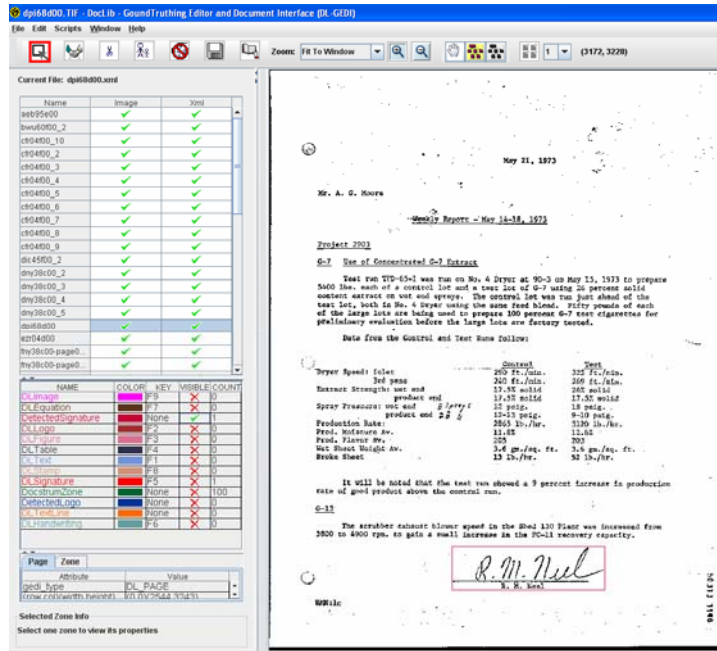
**Figure 2: Visualization of the list of extensible zone types, which include pre-defined zone types in the ground truth and those later defined by the user themselves.**

Figure 3 shows specific zone types displayed using GEDI software.



**(a)**

**(b)**

**Figure 3: Display of zone with select types using GEDI software (a) Page segments. (b) Detected signature regions.**

## 4. Contact Information

For ordering information contact:

```
Laboratory for Language and Media Processing Laboratory
Institute for Advanced Computer Studies
University of Maryland
College Park, Maryland 20742

Attention: Dr. David Doermann

Phone: 301-405-1767
FAX: 443-638-0236
E-mail: doermann@umiacs.umd.edu
http://lamp.cfar.umd.edu
```