# Synergies from Document Analysis

**Daniel P. Lopresti**

`lopresti@cse.lehigh.edu`

Professor, Department of Computer Science and Engineering, Lehigh University

President, International Association for Pattern Recognition (IAPR)

Vice Chair, Computing Research Association's Computing Community Consortium (CCC)

ASQDE_AGM 2.0 ver. 2021 — The Future is Now · August 10-12, 2021

# Goals

- Provide short overview of document analysis research.*
- Survey history of the field briefly.
- Describe some current problems of interest.
- Highlight potential synergies with forensic science.
- Point to online resources that may be useful.
- Offer my thanks to Samiah Ibrahim, ASQDE President, and other conference organizers for providing this opportunity.

\* Cannot completely avoid discussing handwriting, but will try to minimize potential duplication.

# My background

- Professor of Computer Science & Engineering at Lehigh.
- President of International Association for Pattern Recognition.
- Vice Chair of Computing Research Association's CCC Council.
- 30 years in the field; co-EIC of IJDAR; co-Program Chair of ICDAR 2021 and numerous other past conferences.
- B.A. from Dartmouth, Ph.D. from Princeton.
- Research interests include algorithmic and systems-related questions in document analysis, pattern recognition / machine learning, and computer security including electronic voting.

# What is document analysis?

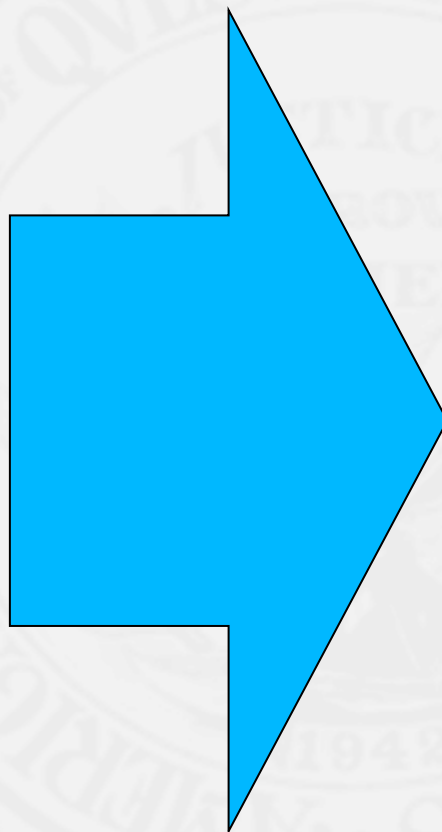Documents are one of humankind's most significant creations:

- Attempt to provide machines with human levels of capability when it comes interpreting documents, broadly defined.

Document analysis combines:

- Computer vision.
- Image processing.
- Machine learning (pattern recognition).
- Natural language understanding.
- Domain expertise.

# What is the target of document analysis?

| | |
|---|---|
| plain text | correct word order for OCR |
| illustrated text | reading order, links to illustrations |
| structured text | compilable or executable form |
| envelope, letter | routing information |
| directory, TOC | name-attribute pairs |
| business form | links to database, add tags |
| schematic diagram | net list or graph |
| engineering drawing | current CAD format |
| map | GIS representation |
| music score | MIDI representation |
| table | layout-independent descriptor |

Adapted from "State of Art of Document Image Processing" (PowerPoint presentation), G. Nagy.

# History of document analysis: a few highlights

- OCR – optical character recognition – dates back to early 1900's.
- Growing practical interest in 1970's (Ray Kurzweil and others).
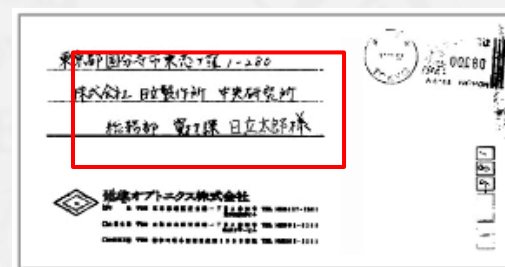
From a
1982 paper:



- Block segmentation
- OCR
- Dictionary check
- Graphics ID
- Vectorization
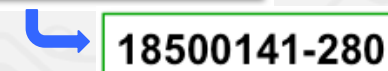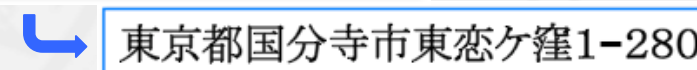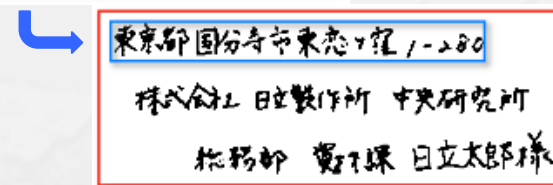- Halftone to grayscale
- Page markup
- Various editors

"Document Analysis System," K. Y. Wong, R. G. Casey, and F. M. Wahl, IBM Journal of Research and Development, vol. 26, no. 6, November 1982, pp. 647-656.

# History of document analysis: a few highlights

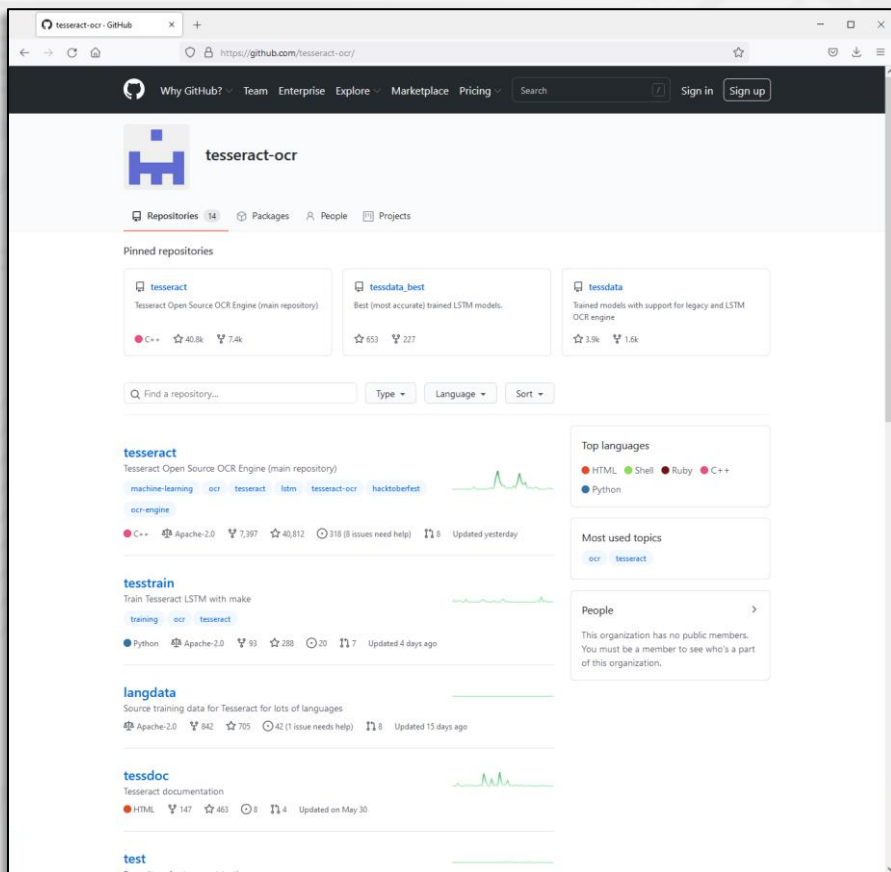Hitachi high-speed postal address reader from 2006:



Hitachi Postal Reader/Sorter
Courtesy Dr. H. Fujisawa

Japanese postal address reading



東京都国分寺市東恋ケ窪1-280

18500141-280

# History of document analysis: a few highlights

Open Source:



Tesseract OCR

# History of document analysis: a few highlights

Seeded Deep Learning revolution:

# More recent topics (from ICDAR 2021 CFP)

- Document image processing
- Text and symbol recognition
- Document analysis systems
- Indexing and retrieval of documents
- Extracting document semantics
- Document summarization and translation
- Human document interaction
- Mobile text recognition
- Scene text detection and recognition
- Recognition of tables and formulas
- Signature verification
- Medical document analysis
- Document analysis for literature search

- Physical and logical layout analysis
- Handwriting recognition
- Document classification
- Document synthesis
- NLP for document understanding
- Office automation
- Multimedia document analysis
- Pen-based document analysis
- Graphics recognition
- Historical document analysis
- Document forensics and provenance analysis
- Document analysis for social good
- Gold-standard benchmarks and data sets

# Is document analysis primarily retrospective?

- Most documents we wish to keep are now produced digitally: books, journals, newspapers, letters, drawings, forms (like tax returns and visa applications).

But …

- Many pre-1980 documents remain to be converted, some of business value (utility drawings), and many historical artifacts.

- Original software or digital medium is not always available (conversion of CAD drawings, tech journals, census data).

- Digital version is not always available $\Rightarrow$ personal applications.

Adapted from "State of Art of Document Image Processing" (PowerPoint presentation), G. Nagy.

# Popular venues: ICDAR

- The International Conference on Document Analysis and Recognition has been held every two years since 1991. Next conference will be held in Lausanne, Switzerland in September.



16th International
Conference on Document
Analysis and Recognition
ICDAR 2021
September 5-10, 2021, Lausanne, Switzerland

- Biggest, broadest conference in document analysis.

- ~200 papers, ~400 attendees.

- Also workshops, competitions, tutorials.

# Popular venues: IWCDF

• ICDAR 2021 Workshop on Computational Document Forensics.



IWCDF2021
in conjunction with ICDAR 2021

Home    Related Topics    Important Dates    Paper Submission    Submit your paper

3rd International Workshop on Computational Document Forensics

05 September 2021, EPFL, Lausanne, Switzerland

Program (soon)

"Everywhere around the world, industries and government processes are being more and more digitized. Document management systems and digital safe-boxes are particularly concerned by these questions, since documents generally remain the basis of many decisions for transactions, contracts and communication. Documents also remain the proofs for many legal issues. As a consequence, it becomes absolutely essential to develop computational forensic science applied to documents and to create the conditions for protecting documents, for confirming their authenticity and for detecting frauds."

# Popular venues: IWCDF

- Prevention of forgeries in documents
- Detection of forged documents
- Detection of fake documents
- Authentication of documents
- Forgery localisation
- Copyright protection
- Watermarking
- Digital signatures
- Taxonomy of features
- Expert results vs. system outputs

- Forensic handwriting verification/identification
- Forensic signature verification/identification
- Within writer versus between writer variations
- Determining the frequency of occurrence of handwriting features
- Automated signature identification and verification
- Automated handwriting identification and verification
- Extraction of movement order features out of the ink trace
- Allograph matching and clustering
- Classification of signatures: legible vs. illegible, complex vs. simple
- Detection of forgeries in printed and rescanned documents

# Popular venue:  DAS

- The Workshop on Document Analysis Systems (DAS) has been held every two years since 1996, most recently virtually in July 2020 (was to have been in Wuhan, China).



- Unlike other more general conferences in the field, DAS focuses on systems-related issues, although often related papers appear as well (e.g., classifier techniques).

# Popular venue: IJDAR



- International Journal on Document Analysis and Recognition (Springer) is devoted to the field.

- Topics include pen-based computing (signature verification) and document authentication / validation, among many others.

# Some recent papers (1)

ICDAR 2019:

- Offline Writer Identification Based on the Path Signature Feature.

- GRK-Papyri: A Dataset of Greek Handwriting on Papyri for the Task of Writer Identification.

- Online Writer Identification using GMM Based Feature Representation and Writer-Specific Weights.

- A Spatio-Spectral Hybrid Convolutional Architecture for Hyperspectral Document Authentication, M. J. Khan, K. Khurshid, F. Shafait.

- Deep Dynamic Time Warping: End-to-End Local Representation Learning for Online Signature Verification.

- Capturing Micro Deformations from Pooling Layers for Offline Signature Verification.

- Offline Signature Verification using Structural Dynamic Time Warping.

- Online Signature Verification by Few-Shot Separable Convolution Based Deep Learning.

# Some recent papers (2)

DAS 2018:

- Encoding CNN Activations for Writer Recognition.

- Gaussian Process Classification as Metric Learning for Forensic Writer Identification.

- Stable Regions and Object Fill-Based Approach for Document Images Watermarking.

- Towards Detection of Morphed Face Images in Electronic Travel Documents, U. Scherhag, C. Rathgeb, and C. Busch.

- A New Descriptor for Pattern Matching: Application to Identity Document Verification.

- Offline Bengali Writer Verification by PDF-CNN and Siamese Net.

- Saliency-Based Detection of Identity Documents Captured by Smartphones.

- Automated Forgery Detection in Multispectral Document Images Using Fuzzy Clustering.

IJDAR March 2020:

- Even big data is not enough: need for a novel reference modelling for forensic document authentication, U. Garain and B. Halder.

# Competitions

- ICDAR 2021 Competition on On-Line Signature Verification.

- ICDAR 2019 Competition on Signature Verification based on an On-line and Off-line Signature Dataset.

- ICDAR 2017 Competition on Historical Document Writer Identification (Historical-WI).

- ICDAR 2017 Competition on Multi-script Writer Identification Using LAMIS-MSHD and CERUG Databases.

- ICDAR 2015 Signature Verification and Writer Identification Competitions for On- and Offline Skilled Forgeries (SigWIComp-2015).

- ICDAR 2015 Multi-Script Writer Identification and Gender Classification (MS-WIGC-2015).

# ICDAR 2017 Historical-WI

- Goal was retrieval of pages which have been written by same author.

- Test dataset consisted of 3,600 handwritten pages originating from 13th to 20th Century.

- Manuscripts from 720 different writers where each writer contributed five pages.

- Five different institutions submitted six methods which were ranked using identification and retrieval metrics.

# Some current problems of interest

As gleaned from upcoming ICDAR 2021 conference:

• Automating document layout analysis.

• Scene text recognition.

• Deep learning approaches for information extraction.

• Making historical manuscript images searchable.

• Extracting data from images of charts / tables / census reports.

• Authorship determination / writer identification.

• Asking questions on document collections / image text.

# Some of my own work over the years (1)

- "Quantifying Information Leakage in Document Redaction," D. Lopresti and A. L. Spitz, Proceedings of the First ACM Workshop on Hardcopy Document Processing (in association [with] and Knowledge Management), November 2004, Washington, DC, p[p. ]

- "Information Leakage Through Document Redaction: Attacks and C[ountermeasures,] Spitz, Proceedings of Document Recognition and Retrieval XII (IS& Electronic Imaging), January 2005, San Jose, CA, pp. 183-190.

- "The Effectiveness of Generative Attacks on an Online Handwriting Biometric," D. Lopresti and J. Raim, Proceedings of the International Conference on Audio- and Video-based Biometric Person Authentication, July 2005, Rye Brook, NY, pp. 1090-1099.

- "Biometric Authentication Revisited: Understanding the Impact of Wo[lf...] Monrose, and D. Lopresti, Proceedings of the Fifteenth USENIX Se[curity...] Vancouver, BC, Canada, pp. 29-41.

- "Evaluating the Security of Handwriting Biometrics," L. Ballard, D. L[opresti...] Tenth International Workshop on Frontiers in Handwriting Recognitio[n...] 461-466.

- "Forgery Quality and its Implications for Behavioral Biometric Security," L. Ballard, D. Lopresti, and F. Monrose, IEEE Transactions on Systems, Man and Cybernetics, Part B, vol. 37, no. 5, October 2007, pp. 1107-1118.

Resources can be combined in an interactive system to undo attempts to hide information via redaction in certain cases

Traditional attack models for behavioral biometric security underestimate talented, resourceful adversaries

# Some of my own work over the years (2)

- "Biometric Key Generation Using Pseudo-Signatures," L. Ballard, J. Chen, D. Lopresti, and F. Monrose, Proceedings of the Eleventh International Conference on Frontiers in Handwriting Recognition, August 2008, Montréal, Canada, pp. 646-651.

- "Pseudo-Signatures as a Biometric," J. Chen, D. Lopresti, L. Ballard, Second IEEE International Conference on Biometrics: Theory, Appl 2008, Arlington, VA, pages 6 (CD-ROM).

> Pseudo-signatures (graphical passwords) can overcome some inherent issues with using real signatures

- "Toward Resisting Forgery Attacks via Pseudo-Signatures," J. Chen, D. Lopresti, and F. Monrose, Proceedings of the Tenth International Conference on Document Analysis and Recognition, July 2009, Barcelona, Spain, pp. 51-55.

- "The Impact of Ruling Lines on Writer Identification," J. Chen, D. Lopresti, and E. Kavallieratou, Proceedings of the Twelfth International Conference on Frontiers in Handwriting Recognition, November 2010, Kolkata, India, pp. 439-444.

- "Using Perturbed Handwriting to Support Writer Identification in the Chen, W. Cheng, and D. Lopresti, Document Recognition and Retri Symposium on Electronic Imaging), January 2011, San Francisco, C

> Methods for improving performance of writer identification in face of severe constraints on training data

- "Parameter Calibration for Synthesizing Realistic-Looking Variability in Offline Handwriting," W. Cheng and D. Lopresti, Document Recognition and Retrieval XVIII (IS&T/SPIE International Symposium on Electronic Imaging), January 2011, San Francisco, CA, pp. 78740Y-1 - 78740Y-10.

# Some of my own work over the years (3)

- "A New Method For Detecting Altered Text in Document Images," L. D. Lopresti, B. Seraogi, and B. B. Chaudhuri, Proceedings of the Se Recognition and Artificial Intelligence (ICPRAI 2020), October 2020

- "Forged Text Detection in Video, Scene, and Document Images," L. S. Raghunandan, U. Pal, T. Lu, and D. Lopresti, to appear in IET Im

> Deep learning combined with targeted image processing steps to detect copy-paste and insertion forgeries

- "Camera-based Ballot Counter," G. Nagy, B. Clifford, A. Berg, G. Saunders, D. Lopresti, and E. Barney Smith, Proceedings of the Tenth International Conference on Document Analysis and Recognition, July 2009, Barcelona, Spain, pp. 151-155.

- "Style-Based Ballot Mark Recognition," P. Xiu, D. Lopresti, H. Baird, of the Tenth International Conference on Document Analysis and Re 216-220.

> Approaches to facilitate the trustworthy capture and reading of election ballots

- "Document Analysis Issues in Reading Optical Scan Ballots," D. Lo, Proceedings of the Ninth IAPR International Workshop on Documer pp. 105-112.

- "Characterizing Challenged Minnesota Ballots," G. Nagy, D. Lopresti, E. H. Barney Smith, and Z. Wu, Document Recognition and Retrieval XVIII (IS&T/SPIE International Symposium on Electronic Imaging), January 2011, San Francisco, CA, pp. 787413-1 - 787413-10.
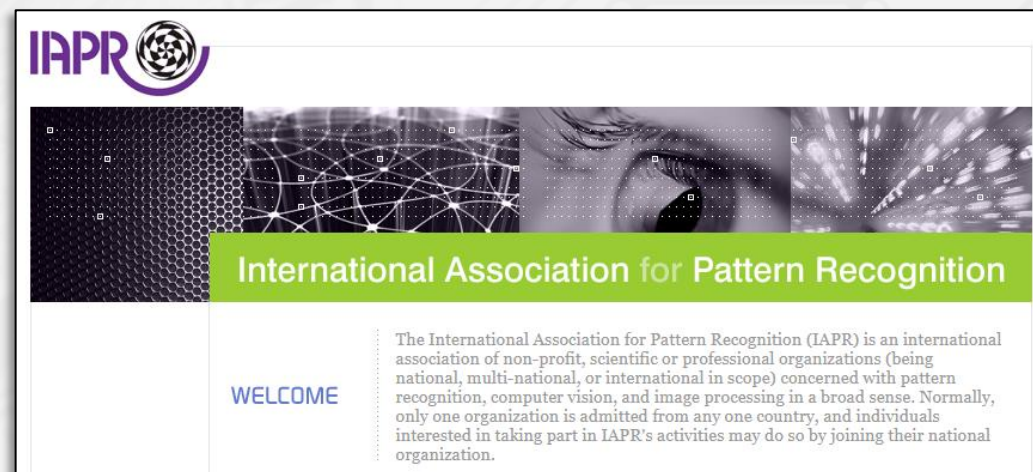
# Other useful resources: IAPR

IAPR – the International Association for Pattern Recognition – is the organization responsible for supporting a wide range of activities by the research community.

• IAPR organizes conferences, administers awards.

IAPR technical committees do much of the actual work:



• TC-06 (Computational Forensics),

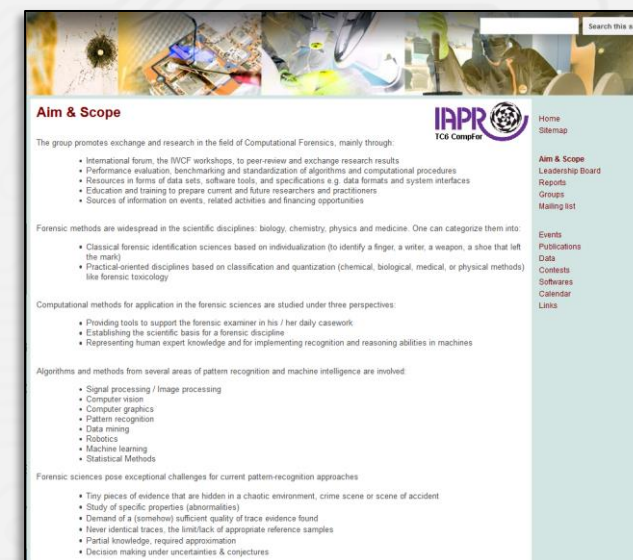• TC-10 (Graphics Recognition),

• TC-11 (Reading Systems).

See www.iapr.org.

# Other useful resources: TC-06

IAPR TC-06 ("Computational Forensics") aims and scope:

- International forum, the IWCF workshops, to peer-review and exchange research results.

- Performance evaluation, benchmarking and standardization of algorithms and computational procedures.

- Resources in forms of data sets, software tools, and specifications e.g. data formats and system interfaces.

- Education and training to prepare current and future researchers and practitioners.

- Sources of information on events, related activities and financing opportunities.

# Other useful resources: TC-10

TC-10 ("Graphics Recognition") promotes interaction among researchers working in document image analysis in general, and graphics recognition in particular.

## Datasets/Softwares
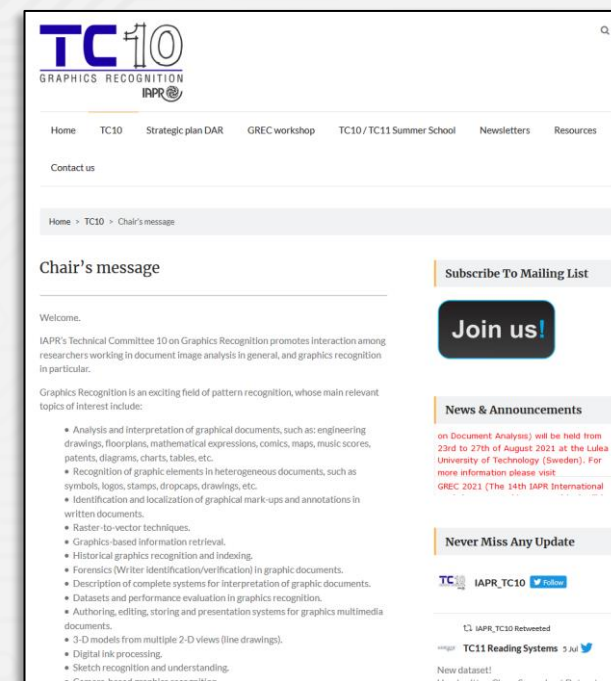
### Engineering drawings / floor plans datasets:
- Bethlehem Steel Dataset (in collaboration with Lehigh University)
- BRIDGE (by Shreya Goyal, Chiranjoy Chattopadhyay) (Paper)
- CVC-FP (Database for structural floor plan analysis)
- FPLAN-POLY
- R-FP-500 (by Rakuten Institute of Technology)
- SESYD (synthetic documents, with the corresponding ground-truth)

### Music Scores datasets:
- List of Music Scores datasets
- ICDAR/GREC competitions on music scores (CVC-MUSCIMA)
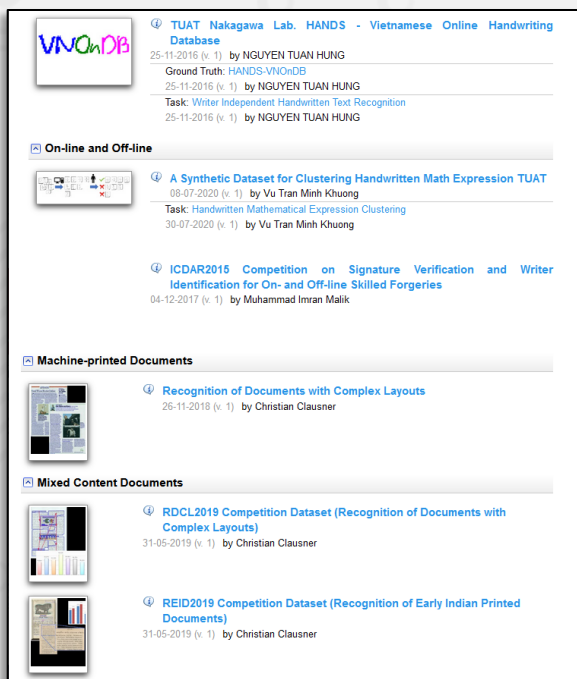
### Comic book datasets:
- BCBID: Bangla Comic Book Image Dataset contains a total of 3327 images of different kinds of 'Bengali Comic Books' from a diverse set of renowned authors.
- COMICS: 1.2 million panels paired with automatic textbox transcriptions from Golden Age collection of the Digital Comics Museum.
- DCM772: 772 annotated images from 27 Golden Age collection of the Digital Comics Museum. It includes ground-truth bounding boxes of all panels, all characters (body + faces), small or big, human-like or animal-like.
- eBDtheque: a representative database of comics of 100 pages including manual annotations of 850 panels and 1092 balloons paired with 1620 comic characters and 4693 text lines.
- FGC 2019 (ICDAR 2019 Competition on Fine-Grained Classification of Comic Characters)
- GNC: the Graphic Narrative Corpus currently contains textual metadata of about 219 titles written in English. Corresponding image are not provided due to copyright issue.
- Manga 109: 109 manga volumes from "Manga Library Z" drawn by professional manga artists in Japan.
- SSGCI 2016 (ICPR 2016 Competition on Subgraph Spotting in Graph representation of Comic Book Images)

Available datasets including engineering drawings, musical scores, very large sets of comic book images.

# Other useful resources: TC-11

TC-11 ("Reading Systems") represents the international research community in topics relating to character recognition and document analysis.





Many available datasets, including signature verification, writer identification, multispectral images of ancient documents.

# Observations and synergies

- Shared interests: shared methods and shared applications.

- Valuable things to be learned looking in both directions.

- Differences in research culture? But these can be bridged.

- Not as much "cross-fertilization" as there could be.

Building connections:

- Sign up for IAPR TC-06 and/or TC-11 mailing lists.

- Peruse resources highlighted in this talk.

- Submit papers to ICDAR (or affiliated workshops), DAS, IJDAR.

- Propose a new, complementary workshop or competition.

# Looking forward

- DAS 2022 will be held in La Rochelle, France in May 2022.

ICDAR 2023

The 17th International Conference on Document Analysis and Recognition

August, 2023 — San Jose, California

- ICDAR 2023 will be held in San Jose, CA in August 2023.
- Calls for Papers, Competitions, and Workshop proposals should be going out soon.

# Thank you

I look forward to the day in the future when we can gather together once again for professional meetings like this.

In the meantime, please feel free to contact me with any questions / comments / suggestions:

lopresti@cse.lehigh.edu

http://www.cse.lehigh.edu/~lopresti/