

CS767 Extra Project

Pragya Sen, BU ID: U97864862

OBJECTIVE: Develop a system that extracts affiliation and contact details of the first and last authors from the enclosed PDF file containing academic paper references and outputs the data into structured Excel files.

DATASET: Used the dataset given with 27 pages worth of references in a PDF file.

METHODOLOGY:

Step 1: PDF OCR Extraction

1. Extracted scanned text using Tesseract OCR.
2. Cleaned text by normalizing quotes and merging broken lines.

Step 2: LLM-Based Reference Parsing

1. Passed reference blocks to LLaMA 3 via Ollama to extract: Reference title, First author, Last author
2. Returned structured JSON for reliability.

Step 3: Multi-Source Online Enrichment

Source	Purpose
Google Scholar	Academic affiliation & verified profile
CrossRef API	DOI inference + affiliation by metadata
DuckDuckGo	Email + affiliation snippets from the web
LLaMA Fallback	Infer missing affiliation/email
Local Cache	Avoid repeated lookups

Step 4: Final Export

Output 2 Excel files:

1. **OCR_output.xlsx** – raw extracted references
2. **Verified_output.xlsx** – enriched author data

Results:

1. OCR_output.xlsx:

	A reference	B title	C first_author	D last_author
2	Wiener '48	Time, communication, and the nervous system. Wiener		
3	[Fuegi '03] Fuegi, J., & Francis, J.	Lovelace & Babbage and the creation of the first computer. Fuegi	Francis	
4	Kagermann '11] Kagermann, H.,	Industrie 4.0: Mit dem Internet der Dinge. Kagermann	Wahlster	
5	Hey '09] Hey, T., Tansley, S., & Tolle, K. M.	The Fourth Paradigm: Data-Intensive Science. Hey, T.	Tolle, K. M.	
6	Mitchell '97	Machine Learning	T.M. Mitchell	
7	[Géron '19] Géron, A. (2019). Hands-on Machine Learning with Python. Géron			
8	[Bayes '63] Bayes, T. (1763). An Essay Towards Solving a Problem in the Doctrine of Chances. Bayes, T.			
9	Russell '21] Russell, S., Norvig, P.	Artificial Intelligence: A Modern Approach. Russell	P. Norvig	
10	[Chapmann '17] Chapman, J.	Machine Learning: Fundamental Algorithms. Chapman	J.	
11	[Knuth '97] Knuth, D. E. (1997).	The Art of Computer Programming. D. E. Knuth		
12	[Bhargava '16] Bhargava, A. (2016). Grokking Algorithms		A. Bhargava	
13	[Han '11] Han, J., Pei, J., & Kamber, M.	Data Mining: Concepts and Techniques. Han	M. Kamber	
14	[Allen '83] Allen, J. F. (1983).	Maintaining knowledge about temporal intervals. Allen	Allen	
15	[Wolpert '97] Wolpert, D.H.	No Free Lunch Theorems for Optimization. Wolpert	Macready	
16	[Topol '19] Topol, E. (2019).	Deep medicine: how artificial intelligence can transform health care. Topol	Topol	
17	[Zheng '18] Zheng, A., & Casari, A.	Feature Engineering for Machine Learning. Zheng, A.	Casari, A.	
18	[Deisenroth '19] Deisenroth, M.P.	Mathematics for Machine Learning. Deisenroth, M.P.	Ong C.S	

2. Verified_output.xlsx:

A reference	B title	C first_author	D last_author	E first_author_affiliation	F first_author_email	G last_author_affiliation	H last_author_email	I
2	Wiener '48 Time, communication, and the nervous system. Wiener			Massachusetts Institute of Technology	Not publicly available (or unknown)			
3	[Fuegi '03] Lovelace & Babbage and the creation of the first computer. Fuegi	Francis		Flare/MIT	fuegi@uwm.edu	Flare/MIT		
4	Kagermann '11] Industrie 4.0: Mit dem Internet der Dinge. Kagermann	Wahlster		Germany	kagermann@acatec.de	Germany	wahler@forschungsgesellschaft-autonomie.de	
5	Hey '09] Hey The Fourth Paradigm: Data-Intensive Science. Hey, T.	Tolle, K. M.		University of Washington		Microsoft Research		
6	Mitchell '97] Machine Learning	T.M. Mitchell				Machine Learning		
7	[Géron '19] Hands-on Machine Learning with Scikit-Learn. Géron							
8	[Bayes '63] An Essay Towards Solving a Problem in the Doctrine of Chances. Bayes, T.							
9	Russell '21] Artificial Intelligence: A Modern Approach. S. Russell	P. Norvig		Stanford University	stuart@cs.stanford.edu	Google Research (formerly pnorvig@norvignon.com)		
10	[Chapmann '17] Machine Learning: Fundamental Algorithms. Chapman	J.		Stanford University (emeritus)	knuth@cs.stanford.edu			
11	[Knuth '97] The Art of Computer Programming. D. E. Knuth							
12	[Bhargava '16] Grokking Algorithms	A. Bhargava		University of California, Berkeley	abhart@eecs.berkeley.edu			
13	[Han '11] Data Mining: Concepts and Techniques	J. Han	M. Kamber	University of Wisconsin-Madison	han@cs.wisc.edu	Monash University	michael.kamber@monash.edu	
14	[Allen '83] Maintaining knowledge about temporal intervals. Allen	Allen						
15	[Wolpert '97] No Free Lunch Theorems for Optimization. Wolpert	Macready		University of Texas at Austin	wolpert@cs.utexas.edu	Sarnoff Corporation	macready@sarnoff.com	
16	[Topol '19] Deep medicine: how artificial intelligence can transform health care. Topol	Topol		Scripps Research Translational Institute	etopol@scripps.edu			
17	[Zheng '18] Feature Engineering for Machine Learning. Zheng, A.	Casari, A.		Carnegie Mellon University	azheng@andrew.cmu.edu	University of Bologna		
18	[Deisenroth '19] Mathematics for Machine Learning	Deisenroth, M.P.	Ong C.S	University of Cambridge	mpd31@cam.ac.uk	National University of Singapore	csong@nus.edu.sg	
19	Saul '00] An Introduction to Locally Linear Embedding. Saul, L. K.	Roweis, S. T.		Massachusetts Institute of Technology	tsaul@mit.edu	University of Texas at Austin	stroewies@utexas.edu	
20	[Van der Maaten '08] Visualizing Data Using t-SNE	Van der Maaten	Hinton	University of Amsterdam	lvmaaten@science.uva.nl	Stanford University	geoffrey.hinton@stanford.edu	
21	[McInnes '18] Umap: Uniform manifold approximation and projection. McInnes	Melville		The University of Melbourne	lachlan.mcinnis@melville.edu.au			