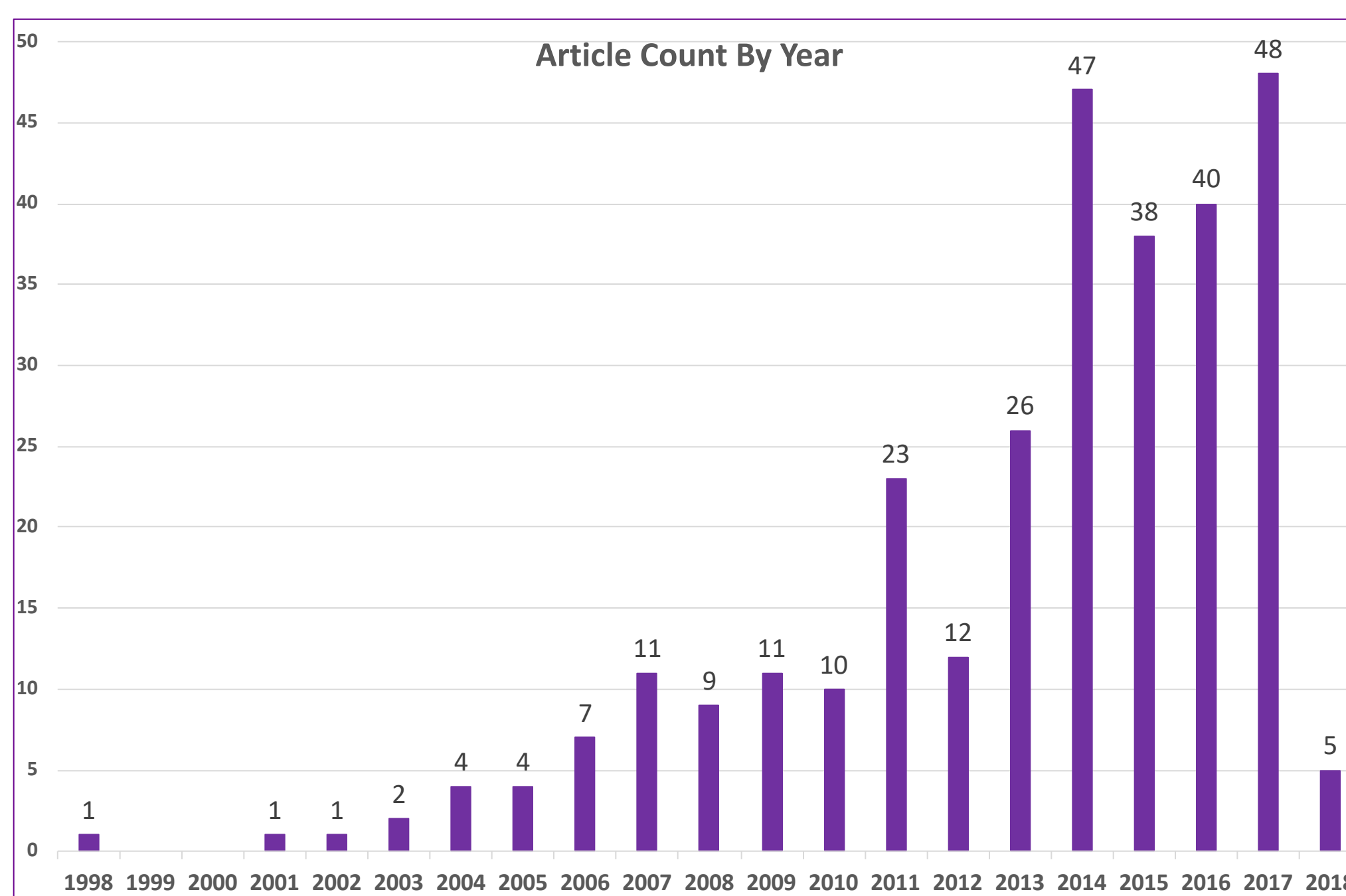




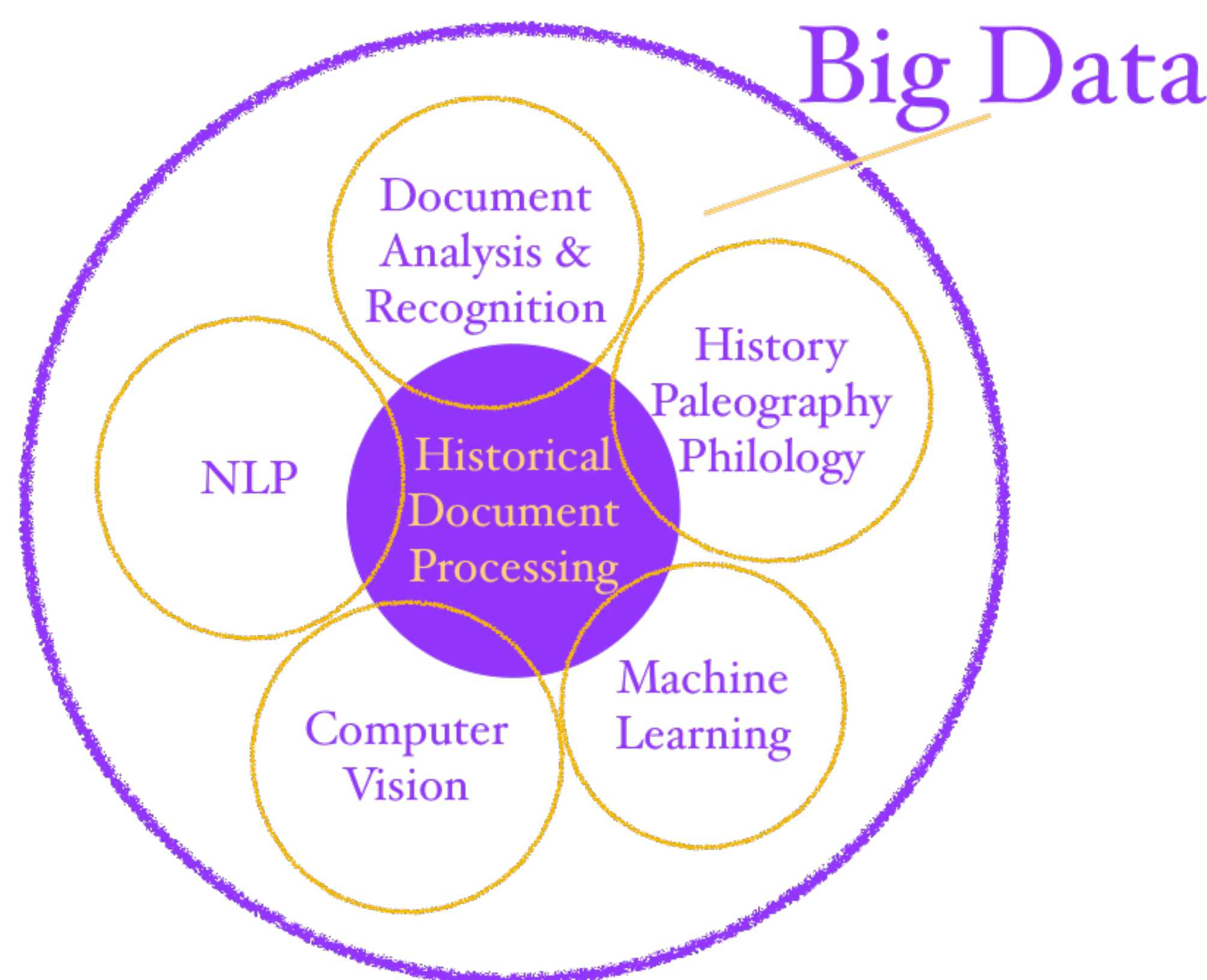
INTRODUCTION

Historical Document Processing is the process of digitizing written material from the past for future use by historians and other scholars. It incorporates algorithms and software tools from various subfields of computer science, including computer vision, document analysis and recognition, natural language processing, and machine learning, to convert images of ancient manuscripts, letters, diaries, and early printed texts automatically into a digital format usable in information retrieval systems. Within the past twenty years, as libraries, museums, and other cultural heritage institutions have scanned an increasing volume of their historical document archives, the need to transcribe the full text from these collections has become acute. Big Data Analytics and infrastructure will be an essential tool in this field. This study compares performance analysis of two OCR systems, discusses HDP workflow, and highlights the role of OCR software in a RESTful API for HDPaaS.

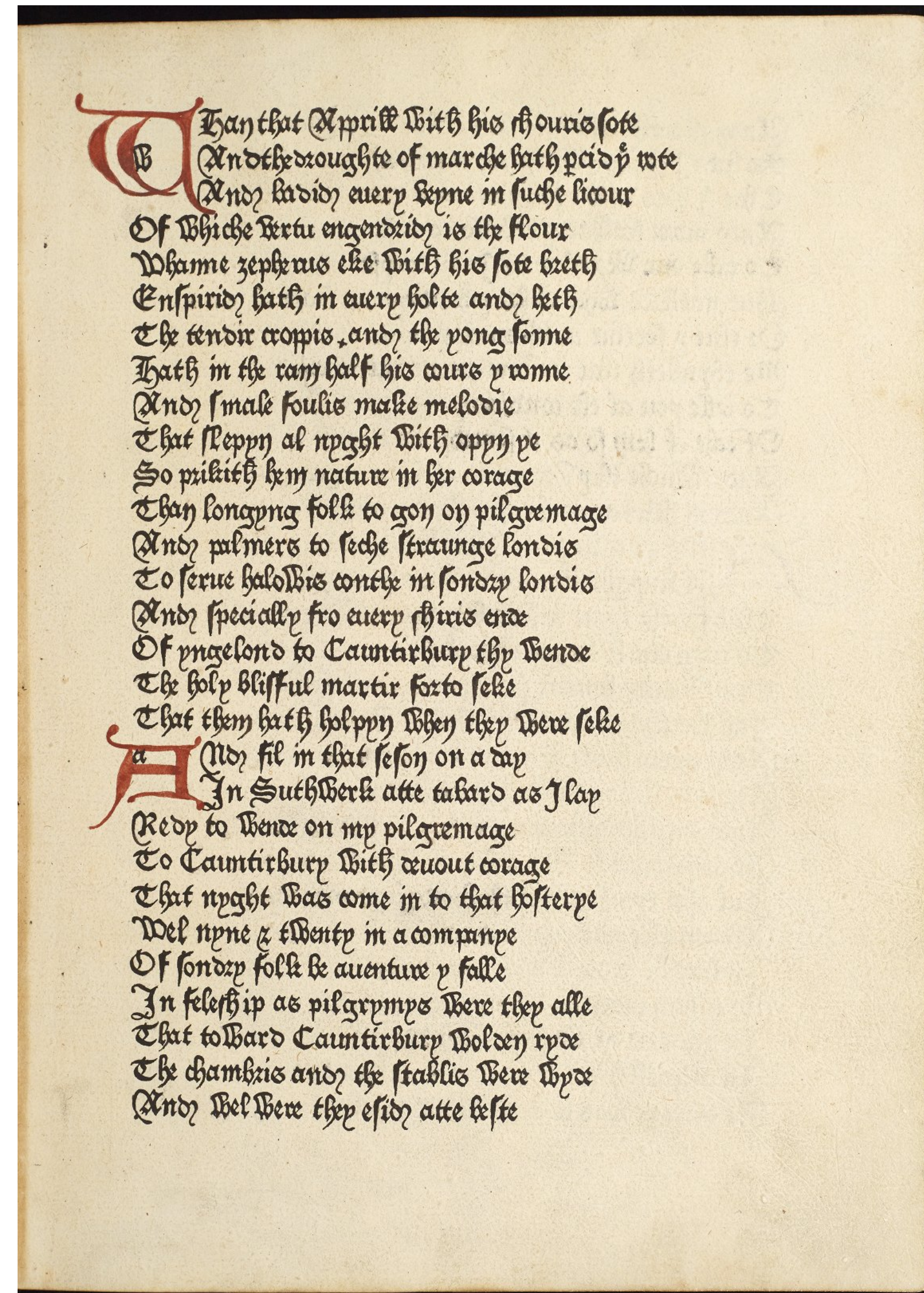
LITERATURE REVIEW



BIG DATA HDP: A HYBRID FIELD



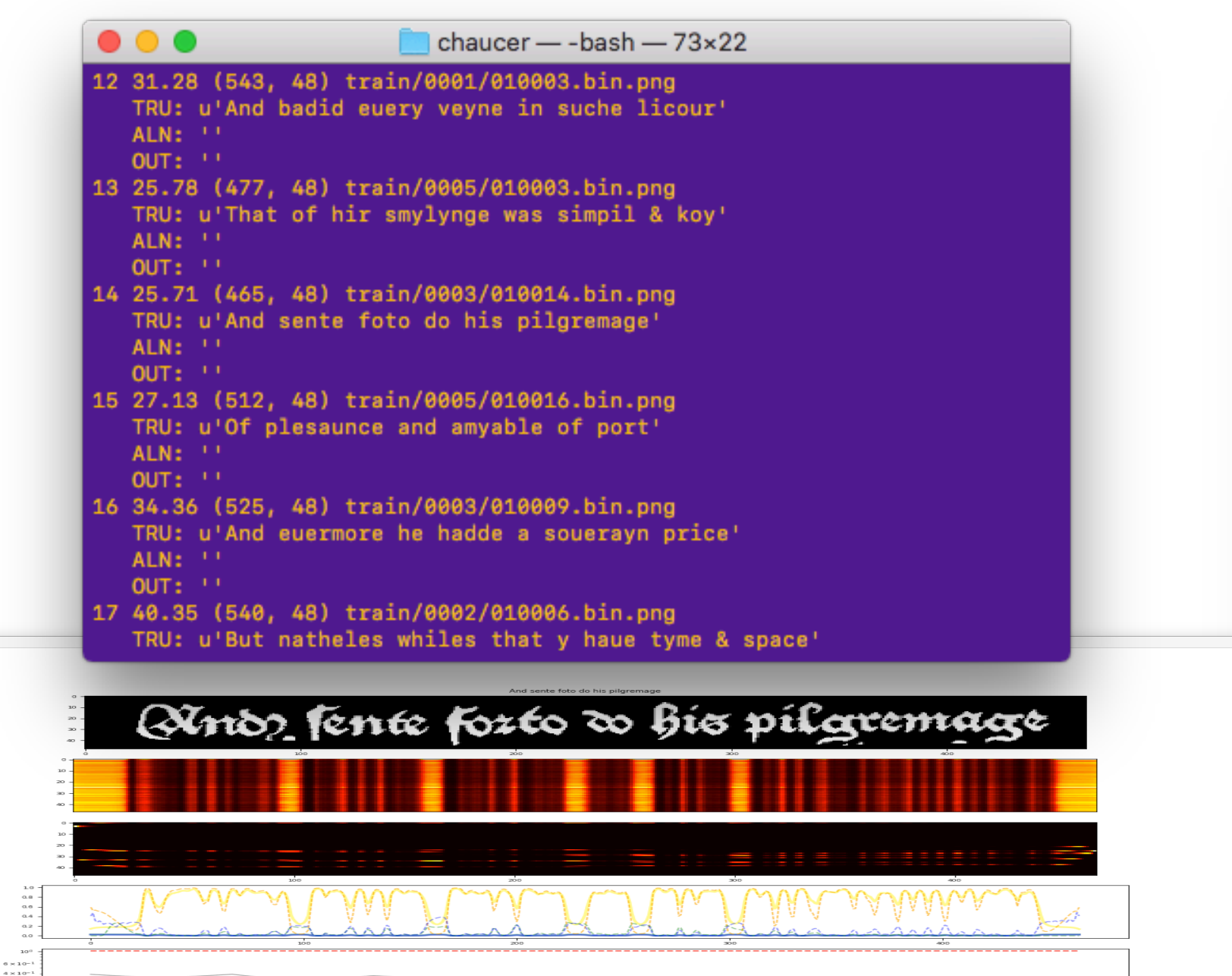
HISTORIC DOCUMENT PROCESSING WORKFLOW



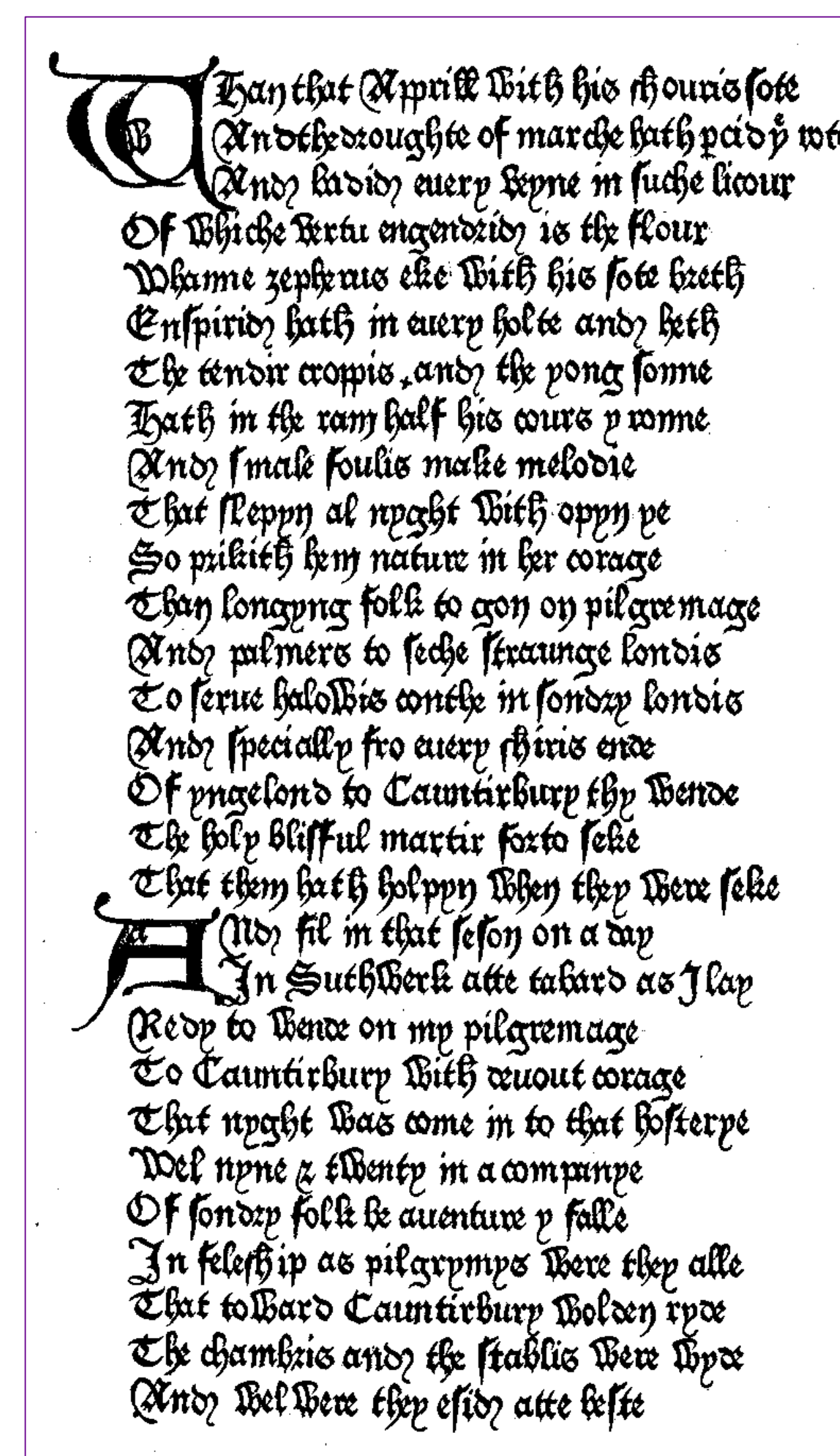
1. ORIGINAL IMAGE

In feleþip as pilgrymys Were they alle
That toþard Cauntirbury Woldeþ ryde
The chambria andþ the stablis Were Wyþe
Andþ Wel Were they eþidþ atte leþte
Andþ kadidþ euery þeyne in ſuche licour
Redy to Wende on my pilgremage
To Cauntirbury With deuout corage

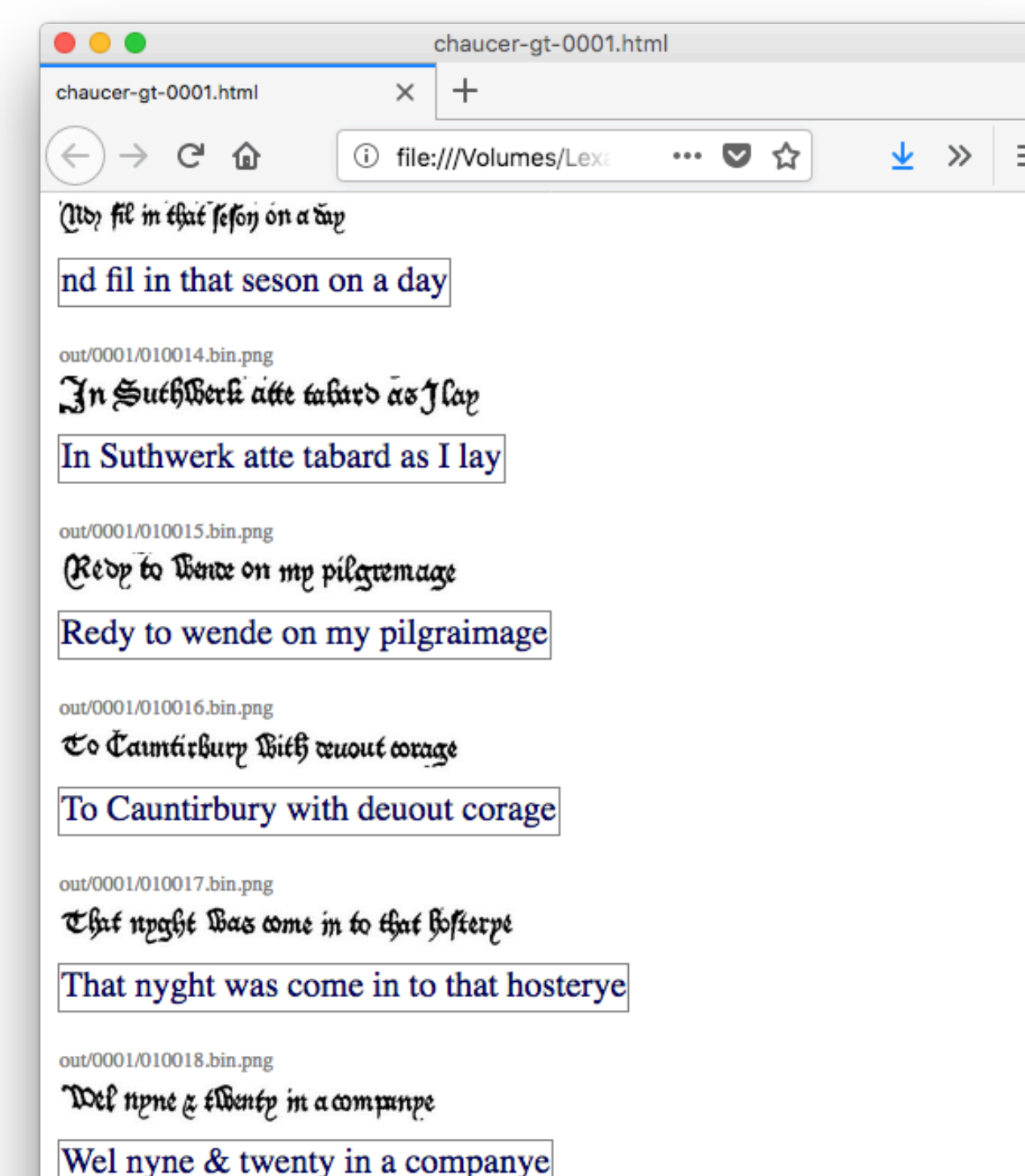
3. LINE SEGMENTATION



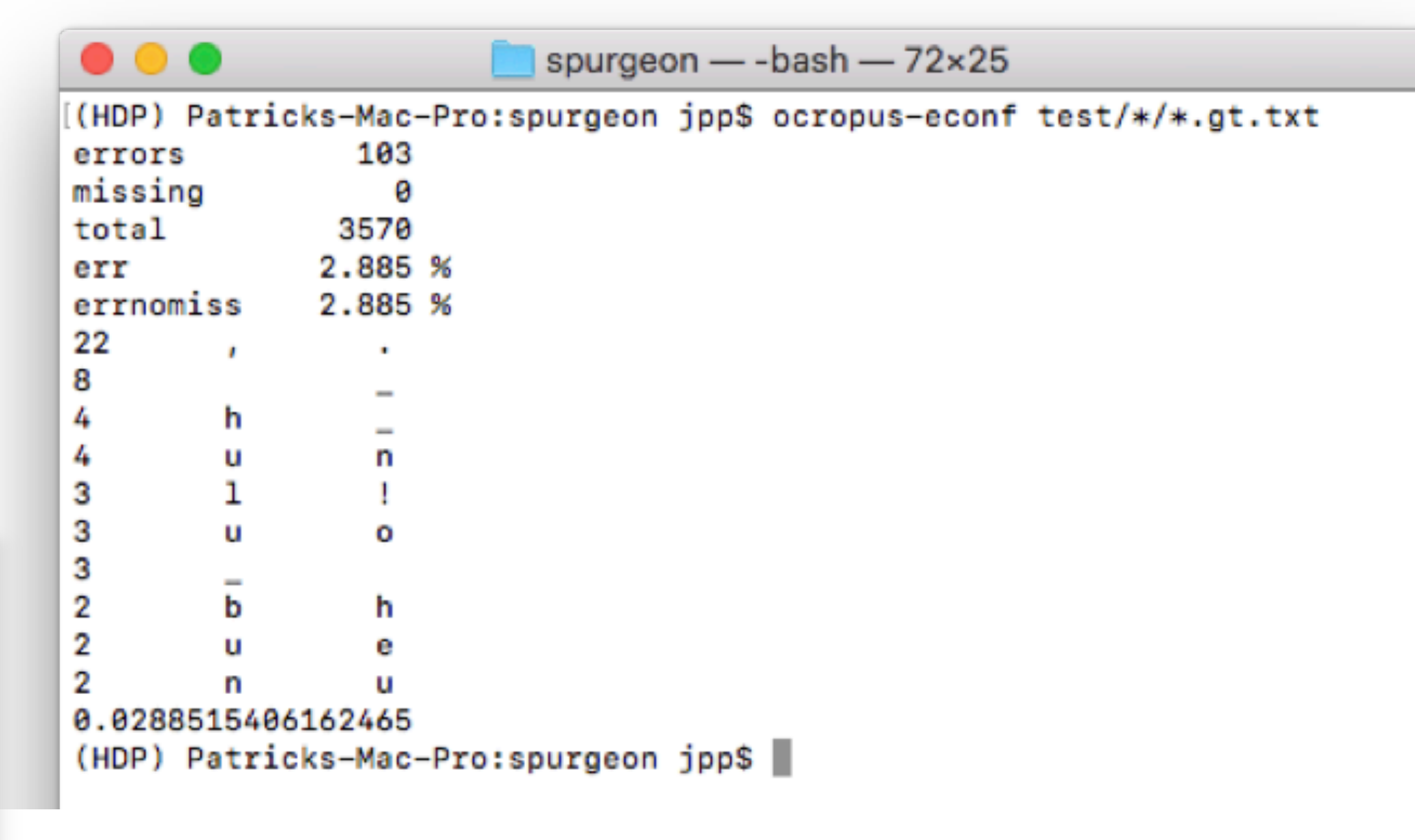
5. TRAINING



2. BINARIZED IMAGE



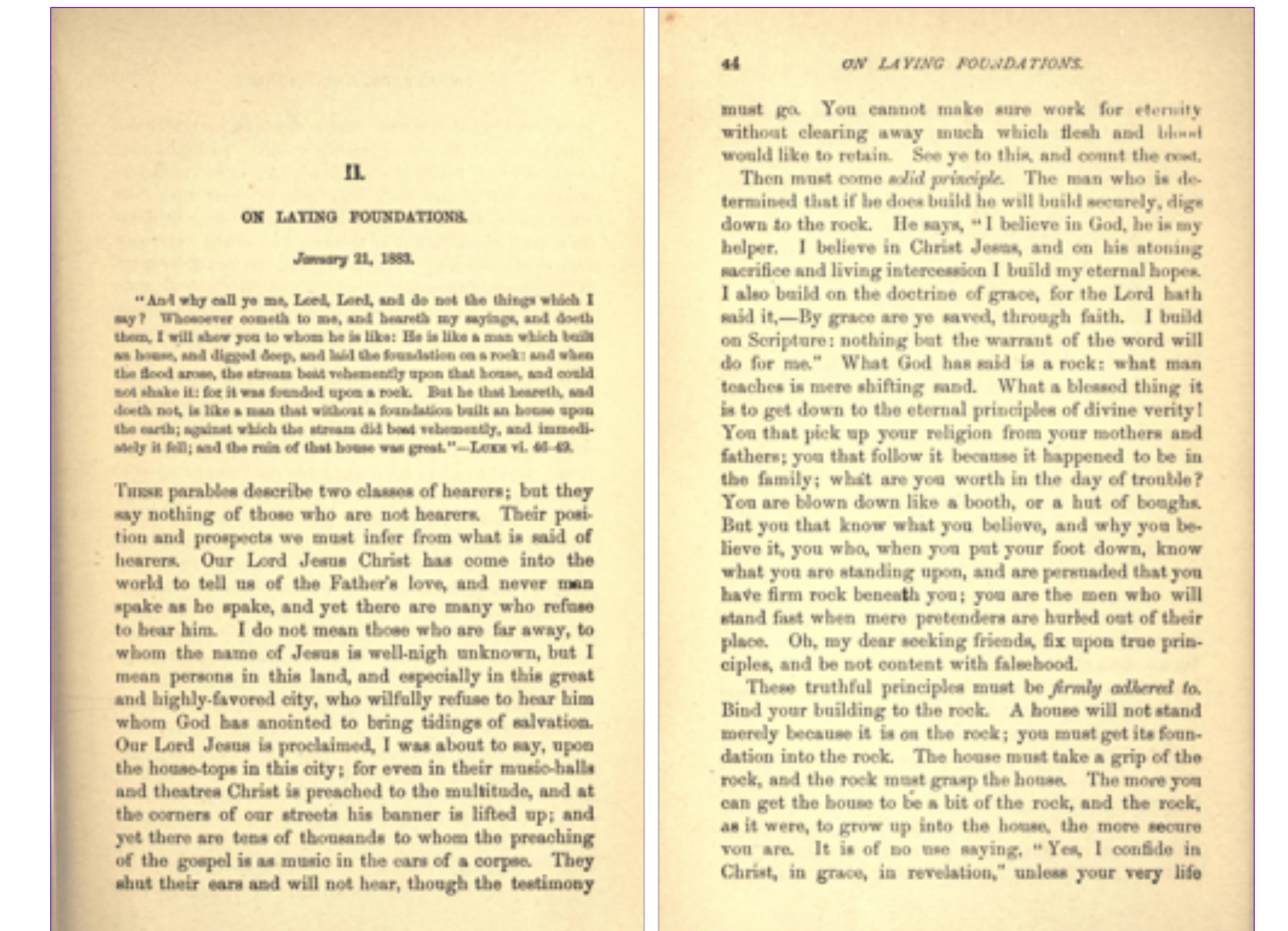
4. GROUND TRUTH ANNOTATION



6. ANALYTICS

CASE STUDY

- Used images from the sermons of C.H. Spurgeon
- Compared OCR performance of OCRopy and Ocular OCR software
- OCRopy uses BLSTM neural network & Ocular uses unsupervised machine learning with a multiple models.



	# Training Lines	# Test Lines	Training Time	CER
OCRopus (trained)	2998*	62	2 hours	33.109%
OCRopus (default model)	NA	3570	NA	2.885%
Ocular	1785	1785	4 hours, 43 minutes**	76.813%

* Some lines duplicated in training data per practitioner recommendations. A model was saved every 1000 lines

** Another training session with a training set 2x larger lasted 17 hours, 33 minutes

- Based on my case study, I have found OCRopy to be a superior OCR system to Ocular due to accuracy and performance metrics
- If OCR software were used in a high performance cluster computing environment, HDP could be implemented with a RESTful API as a cloud service: HDPaaS.
- The extensive quantity of archival data in libraries necessitates a solution using Big Data analytics.

ACKNOWLEDGEMENTS

James Patrick Philips wishes to thank his mentor, Dr. M.H.N. Tabrizi for his inspiration and guidance on this project. He also thanks Rebecca Stamilio-Ehret, Trey Cherry, and Timothy Boyd of Edgecombe Community College for their encouragement and support.