

# Breakdowns in Machine Reading:

---

Attempting to De-privilege Modern English Print with the  
Power of Supercomputing and the DH Dashboard

# The State(s) of the OCR Problem

---

## Computer Science, circa 1990

- The OCR problem is solved! We can machine-read text with 99% accuracy.

## The Humanities, circa 2010

- Wait, what? We still can't reliably perform digital searches on texts printed before 1800! And what about OCR for languages other than modern English?

little better herein, and tell he do it, the opinion of E. P. stands good, for all his pretended self authoritative assertion, and for the close of this Doctors fifth Section, E. P. is not so possessed with a Stoicall Apathy, as not to have a sence of his unsavory Jeare (a Character of the goodness of his cause ) and therefore tells him in his own words in another place, upon another occasion, That he will be silent, rather then imitate that part of his Spirit, though there be in his vindication the advantages of many opposite retortions.

3. Slave, an unfit expression to call the Jewish servant by.

Gattel. Gibieuf. de libert. Dei & creature. l. 1. p. 29. ca. 5. Sect. 5.

The third answer to the Doctors explanation of that place of Exod. 21. 6. That the Jewish Servant did not give up himself to be a slave for ever; but his slavery did terminate in the yeare of Iubite, here the Doctor obstreperously vaunts over his opponent, yet he that takes seriously the weight of this mans arguments, will finde ( I presume ) his virulency and malice poyzing more in the scale, then his reason : I must confes, that if a smooth dialect, and an exquisite texture of words could overcloud a truth, this Doctor can do as much as any towards it : He shall be a slave, &c. it were well the Doctor could afford the servant Jew ( whom the Chaldee renders a son of Israel ) a better terme then slave ( this is a word well becoming them that would not out of necessity, but out of choice, and with all their industry make themselves so ) it may be the

( to ) little better herein, and tell he do it. the opinion of E. P; stands goods-lion-al. his pretended feilds authoritative assertion'- and -for the Jclose Wthijs-Dojëtörjicistb SeHqun; E JI?, not, fo-poffii:ffed Mch. a Sfoizcalliliphthy'dasffiot Jto ii:'ive a fence'of his unsavory Jeare(a Cha-rafler of the goodness of. hisciufe j and therefore tells him in his Jcc own words in an ot.lie r place. upon another occasion' That he will u besilenr' pitcher then imitate rh.-zr. part of his Q. spir?t' though there be J inghis vindication che advantages of many opposite terortions. -The third anfwer. to Jth'e Doctors-explanati'oh of that place of Mve:an 'Exod. 21.--6. That the Iewish. Servant did not give up himstlfto unfit ex- be a stavefor ever;'-but his slavery did terminate in the yeare of flubite, premonto hete rhe Do,flor obstreperously ,vaunts over his jopponen'r,

# The Early Modern OCR Project (eMOP)

---

- Funded by a Mellon grant in 2012
- Train extant, open-source OCR engines to read Early Modern printed documents
- OCR documents from Early English Books Online (EEBO) and Eighteenth Century Collections Online (ECCO), or roughly 45 million page images in two years.



It is incredible, how general a  
Toy Goodman *Fact's* Success crea-  
ted in the City of London; there  
was nothing to be seen or heard the  
next Day, but shaking of Hands,  
Congratulations, Reflections on the  
Danger they had escap'd, and Gra-



To make the wether thinne, and airelike faith Caref  
 unto Where fre, with salt waves meet, and what to  
 betweene oth makes a oord or two III. So farre from  
 owes them in bigger uantities Then they are. Thus  
 and not for hunger a sea Pie Spied through this traite  
 feely where it disputing lay, And tend her doubts a  
 eif, but to the exalters good, As are by great ones,  
 to be the Raifers instrument and food. I. If any kind  
 they neither doe, nor wi Fiers they kill not, nor with  
 nor strive to make a prey f beasts, nor their yong fo  
 pursue not, nor do undertake To spoile the nests in  
 all these unkinde kinds feed upon, To kill them if a  
 fast, lent for their destruion. As if twixt Aire and  
 love, and mens will ever bee. Breake of day. Tis tru  
 thou therefore rise from me Whyould we rise, beca  
 because twas night Love which in spight of darknes  
 des pight of light kee us together. Light hath no tor  
 as well as sie, This were the worst, that it could say

4	a	272	2887	316	2943	0
5	k	320	2887	372	2984	0
6	e	376	2887	424	2943	0
7	t	448	2887	483	2953	0
8	h	488	2887	539	2979	0
9	e	544	2887	590	2943	0
10	w	616	2887	694	2942	0
11	e	696	2887	743	2942	0
12	t	744	2887	775	2951	0
13	h	776	2887	828	2980	0
14	e	832	2887	879	2942	0
15	r	880	2887	915	2941	0
16	t	944	2887	978	2953	0
17	h	984	2887	1035	2978	0
18	i	1040	2887	1062	2976	0
19	n	1064	2887	1117	2944	0
20	n	1120	2887	1173	2944	0
21	e	1176	2887	1224	2943	0
22	,	1224	2862	1251	2907	0
23	a	1280	2887	1328	2945	0
24	n	1328	2887	1386	2943	0
25	d	1392	2887	1451	2982	0
26	a	1480	2886	1529	2944	0
27	i	1536	2886	1563	2979	0
28	r	1568	2886	1611	2945	0
29	e	1616	2886	1664	2942	0



### Results Filter

Ground Truth:  Print Font:

Works:  OCR completed date:

Data Set:  From:  To:

OCR Batch:

### Job Queue

Scheduled: 36  
 In-Progress: 0  
 Await Postprocess: 0  
 Failed: 24351  
 Ingested: 40842453  
 Ingest Failed: 0  
**TOTAL:** 40866840

Search:  Show  entries

	Status		Data Set	ID	TCP Number	Title	Author	Font	OCR Date	OCR Engine	OCR Batch	Juxta	RETAS
<input type="checkbox"/>	0-0-63-0		EEBO	1		Alkinoou eis ta tou Platonos eisagoge Alcinoi in Platonicam philosophiam introductio.	Albinus.		03/13/2015 17:30	Tesseract	4: EEBO w/o GT (SC8b-R8-D2b)	N/A	N/A
<input type="checkbox"/>	0-0-2-0		EEBO	2	A24326	An Account of a horrid and barbarous murder committed on the body of a young person supposed to be of a good quality in the fields beyond Whitechappel-Church in the Parish of Stepny ...	Anon.	all other EEBO w/GT	03/11/2015 16:04	Tesseract	3: EEBO with GT (SC8b-R8-D2b)	0.694	N/A
<input type="checkbox"/>	0-0-137-0		EEBO	3		The way of the Spirit in bringing souls to Christ set forth in X sermons on John 16:7, 8, 9, 10 and chap 7:37 / by Mr. Thomas Allen, late pastor of a church in ... Norwich.	Allen, Thomas, 1608-1673.		03/13/2015 17:30	Tesseract	4: EEBO w/o GT (SC8b-R8-D2b)	N/A	N/A
<input type="checkbox"/>	0-0-1-0		EEBO	4	A47215	A contemplation on Bassets-down-Hill by the most sacred adorer of the Muses, Mrs. A.K.	A. K., Mrs.	all other EEBO w/GT	03/11/2015 16:04	Tesseract	3: EEBO with GT (SC8b-R8-D2b)	0.046	N/A
<input type="checkbox"/>	0-0-18-0		EEBO	5	A29730	A dissuasive from popery sent in a letter from A.B. to C.D.	A. B.	all other EEBO w/GT	03/11/2015 16:04	Tesseract	3: EEBO with GT (SC8b-R8-D2b)	0.752	N/A
<input type="checkbox"/>	0-0-3-0		EEBO	6	A24571	An Account of Oliver Hawley and John Condon who were executed at Tyburn on Friday the 2d of July 1686 for robing His Majesties male near Ilford in the county of Essex.	Anon.	all other EEBO w/GT	03/11/2015 16:04	Tesseract	3: EEBO with GT (SC8b-R8-D2b)	0.279	N/A



# eMOP Project OCR Accuracy

---

EEBO

**68%**

ECCO

**86%**

# eMOP Assets ([emop.tamu.edu](http://emop.tamu.edu))

## Database Schema

A MySQL database initially designed for EEBO and ECCO document metadata, then expanded to include ground truth and OCR results.

## Training Data

High-res, cleaned manuscript images with each character meticulously delineated and identified for a multitude of fonts used by Early Modern printers.

## Storage

Network Attached Storage designed to serve up many terabytes of data.

## eMOP Dashboard

A Ruby-on-Rails web app for queueing up OCR jobs and viewing results.

## eMOP Controller

A Python software package for the parallelized running of tasks on the Brazos Supercomputing Cluster.

## OCR Workflow

The establishment of a step-by-step, automated process for OCR'ing page images and then improving that OCR via post-processing.

# Enter the PrimerosLibros project...

---

The IDHMC's first attempt to leverage eMOP assets toward a new endeavor 😊

# Responding to PrimerosLibros Challenges

---

## **Machine-reading in languages other than English (Spanish, Latin, Nahuatl, Huastec, Mixtec, Otomi, Tarascan, and Zapotec)**

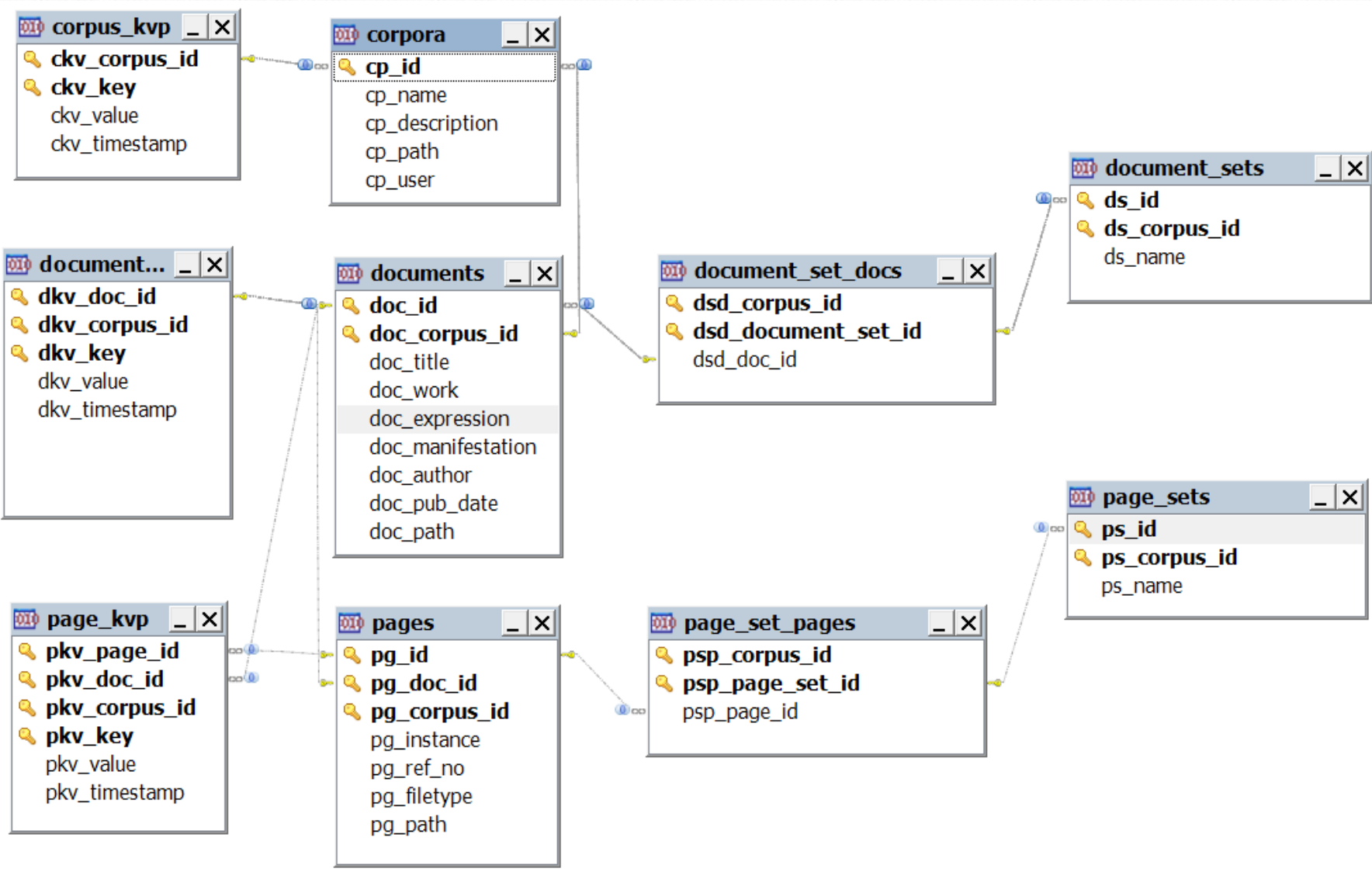
No more code written in older versions of Python (2.7) that did not readily support Unicode out of the box. All tools developed for the PrimerosLibros project use Python 3.5.

# Responding to PrimerosLibros Challenges

---

## **Supporting a New Database Schema**

A complete redesign of the original eMOP database which supports (but does not enforce) a FRBR (work, expression, manifestation, instance) hierarchical structure, and otherwise allows for the storing of key/value pairs at the corpus, document, and page level, providing for a “schemaless” approach.



## All Documents

Show 25 out of 135 rows

	Document ID ▲	FB ID ▲	Title ▲	Printer ▲	Year ▲	Training Set ▲	View ▲
<input type="checkbox"/>	1000	pl_blac_004	<i>Doctrina cristiana para instrucción e informaci...</i>	J. Cromberger	1544	93	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1113	pl_blac_001	<i>Doctrina breve muy provechosa de las cosas q...</i>	J. Cromberger	1544	91	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1114	pl_blac_002	<i>Tripartito del cristianísimo y consolatorio docto...</i>	J. Cromberger	1544	158	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1115	pl_blac_003	<i>Este es un compendio breve que trata de la ma...</i>	J. Cromberger	1544	24	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1124	pl_houg_001	<i>Doctrina breve muy provechosa de las cosas q...</i>	J. Cromberger	1544	169	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	893	pl_jcbl_002	<i>Cartilla para la enseñanza de la doctrina cristia...</i>	J. Pablos	1547	undefined	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1020	pl_jcbl_001	<i>Doctrina cristiana en lengua mexicana</i>	J. Cromberger /...	1547	40	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1252	pl_blac_020	<i>Regla cristiana breve para ordenar la vida y tie...</i>	undefined	1547	94	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1250	pl_blac_018	<i>Doctrina cristiana en lengua española y mexica...</i>	J. Pablos	1550	21	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	1251	pl_blac_019	<i>Doctrina cristiana en lengua mexicana</i>	J. Pablos	1553	22	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	901	pl_blac_005	<i>Recognitio, summularum reverendi patris Illdep...</i>	J. Pablos	1554	25	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	902	pl_blac_006	<i>Dialectica resolutio cum textu Aristotelis</i>	J. Pablos	1554	26	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	959	pl_uslp_001	<i>Recognitio, summularum reverendi patris Illdep...</i>	J. Pablos	1554	86	<a href="#">Pages</a> <a href="#">Details</a>

## Documents from Set 3 Engine Comparison

[Back to All Documents](#)Show  out of 61 rows

<input type="checkbox"/>	Doc ID	Title	GCV Score	Tess4 Score	eMOP Score	Source	View
<input type="checkbox"/>	40355	<i>Pasquils passe, and passeth not Set downe in three pees. His passe, precession, and progn...</i>	0.25	0.30	0.78	EEBO	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	40356	<i>Much adoe about nothing As it hath been sundrie times publikely acted by the right honoura...</i>	0.50	0.29	0.81	EEBO	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	40357	<i>The lettin[g] of humours blood in the head-vaine with a new morissco, daunced by seauen s...</i>	0.25	0.21	0.22	EEBO	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	40363	<i>Inimicus amicus an excellent treatise, shewing, how a man may reape profit by his enemy.</i>	0.33	0.22	0.64	EEBO	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	40364	<i>The exaltation of the kingdome and priesthood of Christ In certaine sermons vpon the 110. P...</i>	0.42	0.21	0.57	EEBO	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	40366	<i>A treatise concerning the trinitie of persons in vnitie of the deitie Written to Thomas Manneri...</i>	0.21	0.24	0.68	EEBO	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	40367	<i>Dialogicall discourses of spirits and diuels declaring their proper essence, natures, dispositi...</i>	0.40	0.20	0.66	EEBO	<a href="#">Pages</a> <a href="#">Details</a>
<input type="checkbox"/>	40368	<i>The mirror of diuine prouidence Containing a collection of Theodoret his arguments: declari...</i>	0.08	0.19	0.30	EEBO	<a href="#">Pages</a> <a href="#">Details</a>



# FirstBooks

## FirstBooks/PrimerosLibros

Edit Key-Value Pairs  
Run Corpus Task

### All Page Sets

Show 25 out of 144 rows

Set ID	Name
<input type="checkbox"/> 52	pl_bjml_001
<input type="checkbox"/> 53	pl_bjml_002
<input type="checkbox"/> 150	pl_bjml_003
<input type="checkbox"/> 54	pl_bjml_003
<input type="checkbox"/> 157	pl_bjml_003
<input type="checkbox"/> 90	pl_bjml_004
<input type="checkbox"/> 153	pl_bjml_004
<input type="checkbox"/> 55	pl_bjml_005_TrainingSet_7.5
<input type="checkbox"/> 154	pl_bjml_005_TrainingSet_8.15

### Corpus Key-Value Pairs

Key

corpus\_manager\_tables

Value

```
117 },
118 "documents": {
119   "table_name": "documents",
120   "title": "All Documents",
121   "columns": {
122     "05": {
123       "template": "{{=it.kvp.printer}}",
124       "title": "Printer",
125       "sorter": "string",
126       "sortable": true,
127       "align": "left",
128       "width": 125,
129       "field": "printer"
130     },
131     "04": {
132       "template": "<b><i>{{=it.title}}</i></b>",
133       "title": "Title",
134       "sorter": "string",
135       "sortable": true,
```

Cancel Update

# Responding to PrimerosLibros Challenges

---

## **Supporting Ocular (cont.)**

Because Ocular is under active development, PrimerosLibros collaborators needed a way to pull the latest code down from their git repo, compile, and then install the new version for use.

# FirstBooks

## FirstBooks/PrimerosLibros

Edit Key-Value Pairs  
Run Corpus Task

### All Page Sets

Show 25 out of 144 rows

	Set ID	Name
<input type="checkbox"/>	52	pl_bjml_001_trainingSet_7.5
<input type="checkbox"/>	53	pl_bjml_002_TrainingSet_7.5
<input type="checkbox"/>	150	pl_bjml_003-line-extractor-test-8.4
<input type="checkbox"/>	54	pl_bjml_003_TrainingSet_7.5
<input type="checkbox"/>	157	pl_bjml_003_TrainingSet_8.21
<input type="checkbox"/>	90	pl_bjml_004-TrainingSet_7.11
<input type="checkbox"/>	153	pl_bjml_004-TrainingSet_8.20

### Run a Task

**Job Name**

**Task**

**Job Site**

# Responding to PrimerosLibros Challenges

---

## Supporting a New OCR Engine (Ocular)

Unlike Tesseract, which requires a tedious, manual training process involving pre-processed images, etc., Ocular allows for the unsupervised training of its classifier given a set of exemplary images for each document.

## All Page Sets

Show 25 out of 144 rows

	Set ID	Name
<input type="checkbox"/>	52	pl_bjml_001_trainingSet_7.5
<input type="checkbox"/>	53	pl_bjml_002_TrainingSet_7.5
<input type="checkbox"/>	150	pl_bjml_003-line-extractor-test-8.4
<input type="checkbox"/>	54	pl_bjml_003_TrainingSet_7.5
<input type="checkbox"/>	157	pl_bjml_003_TrainingSet_8.21
<input type="checkbox"/>	90	pl_bjml_004-TrainingSet_7.11
<input type="checkbox"/>	153	pl_bjml_004-TrainingSet_8.20
<input type="checkbox"/>	55	pl_bjml_005_TrainingSet_7.5
<input type="checkbox"/>	151	pl_bjml_005_TrainingSet_8.15
<input type="checkbox"/>	154	pl_bjml_005_TrainingSet_8.20
<input type="checkbox"/>	56	pl_bjml_006_TrainingSet_7.5
<input type="checkbox"/>	165	pl_bjml_006_TrainingSet_8.23
<input type="checkbox"/>	91	pl_blac_001-TrainingSet_7.11
<input type="checkbox"/>	92	pl_blac_002-TrainingSet_7.11

# Responding to PrimerosLibros Challenges

---

## **Supporting Ocular (cont.)**

Whereas the old eMOP Dashboard supported a single, fixed workflow for OCR'ing documents, Ocular requires a two-step process for reading each document (font training on exemplary page sets, transcription).

Show 25 out of 135 rows

Run a Task

Job Name

Create training set for Document 1020

Job Site

Brazos Supercomputing Cluster

Task

- ✓ Create Ocular Training Set
- OCR Document with Tesseract
- OCR Document with Ocular
- OCR Document with Tesseract 4
- Tesseract 4 Post-Processing
- OCR Doc w Tess4 and PostProc
- OCR Document with Google Cloud Vision
- Perform Juxta Comparison

Training Set

View

Training Set	View
93	<a href="#">Pages</a> <a href="#">Details</a>
91	<a href="#">Pages</a> <a href="#">Details</a>
158	<a href="#">Pages</a> <a href="#">Details</a>
24	<a href="#">Pages</a> <a href="#">Details</a>
169	<a href="#">Pages</a> <a href="#">Details</a>
undefined	<a href="#">Pages</a> <a href="#">Details</a>
40	<a href="#">Pages</a> <a href="#">Details</a>
94	<a href="#">Pages</a> <a href="#">Details</a>
21	<a href="#">Pages</a> <a href="#">Details</a>
22	<a href="#">Pages</a> <a href="#">Details</a>
25	<a href="#">Pages</a> <a href="#">Details</a>
26	<a href="#">Pages</a> <a href="#">Details</a>
86	<a href="#">Pages</a> <a href="#">Details</a>
87	<a href="#">Pages</a> <a href="#">Details</a>
39	<a href="#">Pages</a> <a href="#">Details</a>
159	<a href="#">Pages</a> <a href="#">Details</a>

	Document ID	FB ID			
<input type="checkbox"/>	1000	pl_blac_004			
<input type="checkbox"/>	1113	pl_blac_001			
<input type="checkbox"/>	1114	pl_blac_002			
<input type="checkbox"/>	1115	pl_blac_003			
<input type="checkbox"/>	1124	pl_houg_001			
<input type="checkbox"/>	893	pl_jcbl_002			
<input checked="" type="checkbox"/>	1020	pl_jcbl_001			
<input type="checkbox"/>	1252	pl_blac_020			
<input type="checkbox"/>	1250	pl_blac_018			
<input type="checkbox"/>	1251	pl_blac_019			
<input type="checkbox"/>	901	pl_blac_005			
<input type="checkbox"/>	902	pl_blac_006	<i>Dialectica resolutio cum textu Aristotelis</i>	J. Pablos	1554
<input type="checkbox"/>	959	pl_uslp_001	<i>Recognitio, summularum reverendi patris Illdep...</i>	J. Pablos	1554
<input type="checkbox"/>	960	pl_uslp_002	<i>Dialectica resolutio cum textu Aristotelis</i>	J. Pablos	1554
<input type="checkbox"/>	1255	pl_blac_023	<i>Commentaria in Ludovici Vives Exercitationes li...</i>	J. Pablos	1554
<input type="checkbox"/>	1001	pl_blac_007	<i>Aquí comienza un vocabulario en la lengua cas...</i>	J. Pablos	1555

## All Jobs

 Show  out of 339 rows

<input type="checkbox"/>	ID	Job Name	Task Name	Target	Submitted	Status	View
<input type="checkbox"/>	1542	pl_usal_005-ocr-...	OCR Document ...	Document: Vida ...	09-04-2017 11:3...	Results processed (09-07-2017 18:45:55)	<a href="#">Logs</a>
<input type="checkbox"/>	1541	pl_uiab_011-ocr-...	OCR Document ...	Document: Prem...	09-04-2017 11:3...	Results processed (09-07-2017 00:29:53)	<a href="#">Logs</a>
<input type="checkbox"/>	1539	pl_houg_002-ocr-...	OCR Document ...	Document: Provi...	09-04-2017 11:3...	Results processed (09-07-2017 18:48:46)	<a href="#">Logs</a>
<input type="checkbox"/>	1536	pl_bjml_002-ocr-...	OCR Document ...	Document: Tabul...	09-04-2017 10:4...	Results processed (09-05-2017 04:30:11)	<a href="#">Logs</a>
<input type="checkbox"/>	1535	pl_plfx_007-ocr-...	OCR Document ...	Document: Trata...	09-04-2017 10:4...	Results processed (09-05-2017 04:31:42)	<a href="#">Logs</a>
<input type="checkbox"/>	1531	pl_blac_039-ocr-...	OCR Document ...	Document: Psal...	09-03-2017 11:0...	Results processed (09-04-2017 06:50:48)	<a href="#">Logs</a>
<input type="checkbox"/>	1530	pl_blac_032-ocr-...	OCR Document ...	Document: Arte ...	09-03-2017 11:0...	Results processed (09-04-2017 01:47:32)	<a href="#">Logs</a>
<input type="checkbox"/>	1529	pl_blac_023-ocr-...	OCR Document ...	Document: Com...	09-03-2017 11:0...	Results processed (09-04-2017 00:14:18)	<a href="#">Logs</a>
<input type="checkbox"/>	1528	pl_blac_022-ocr-...	OCR Document ...	Document: Manu...	09-03-2017 11:0...	Results processed (09-03-2017 16:12:25)	<a href="#">Logs</a>
<input type="checkbox"/>	1527	pl_blac_016-ocr-...	OCR Document ...	Document: Arte ...	09-03-2017 11:0...	Results processed (09-03-2017 12:39:07)	<a href="#">Logs</a>
<input type="checkbox"/>	1526	pl_uslp_002-train...	Create Ocular Tr...	Document: Diale...	09-02-2017 01:3...	Results processed (09-02-2017 09:07:32)	<a href="#">Logs</a>
<input type="checkbox"/>	1525	pl_uslp_001-train...	Create Ocular Tr...	Document: Reco...	09-02-2017 01:3...	Results processed (09-02-2017 05:58:32)	<a href="#">Logs</a>
<input type="checkbox"/>	1524	pl_usal_005-train...	Create Ocular Tr...	Document: Vida ...	09-02-2017 01:3...	Results processed (09-02-2017 15:13:32)	<a href="#">Logs</a>



Show 25 out of 339 rows

ID	Job Name	Task Name	View
1542	pl_usal_005-ocr...	OCR D	Logs
1541	pl_uiab_011-ocr...	OCR D	Logs
1539	pl_houg_002-ocr...	OCR D	Logs
1536	pl_bjml_002-ocr...	OCR D	Logs
1535	pl_plfx_007-ocr...	OCR D	Logs
1531	pl_blac_039-ocr...	OCR D	Logs
1530	pl_blac_032-ocr...	OCR D	Logs
1529	pl_blac_023-ocr...	OCR D	Logs
1528	pl_blac_022-ocr...	OCR D	Logs
1527	pl_blac_016-ocr...	OCR D	Logs
1526	pl_uslp_002-train...	Create	Logs
1525	pl_uslp_001-train...	Create	Logs
1524	pl_usal_005-train...	Create	Logs
1523	pl_uiab_011-train...	Create	Logs
1522	pl_tlal_008-train...	Create	Logs
1521	pl_tlal_007-train...	Create Ocular Tr... Document: Voca... 09-01-2017 17:3... Results processed (09-01-2017 23:35:04)	Logs
1520	pl_tlal_006-train...	Create Ocular Tr... Document: Voca... 09-01-2017 17:2... Results processed (09-01-2017 23:34:32)	Logs

## Job Configuration

Name

pl\_uiab\_011-train-9.1.17-mvf

```
55 - "parameters": {
56 -   "document_pageset_key": {
57 -     "type": "static",
58 -     "value": "training_set"
59 -   },
60 -   "allow_overlapping_lines": {
61 -     "key": "allow_overlapping_lines",
62 -     "type": "choice_from_kv_list",
63 -     "kv_path": "target.corpus",
64 -     "label": "Allow Overlapping Lines",
65 -     "value": "false"
66 -   },
67 -   "continue_from_last_complete_iteration": {
68 -     "type": "static",
69 -     "value": "true"
70 -   },
71 -   "input_language_model_path": {
72 -     "key": "ocular_lms",
73 -     "type": "choice_from_kv_list",
```

Cancel

Update

Page Num: 12

ID: 101890 Image: pl\_blac\_019\_00012-1000.jpg

fb\_page\_id: pl\_blac\_019\_00012

pl\_blac\_019-ocr-7.19.17-  
mvf(1325)

1000\_transcription\_normalized.txt

**Ocular Transcription:** pl\_blac\_019\_00011-  
1000\_transcription.txt

**Ocular Comparison:** pl\_blac\_019\_00012-  
1000\_comparisons.txt

**Ocular DIPL XML:** pl\_blac\_019\_00012-1000\_dipl.alto.xml

**Ocular NORM XML:** pl\_blac\_019\_00012-1000\_norm.alto.xml

**Ocular Normalized Transcription:** pl\_blac\_019\_00012-  
1000\_transcription\_normalized.txt

**Ocular Transcription:** pl\_blac\_019\_00012-  
1000\_transcription.txt

pl\_blac\_019-OCR\_5.26-  
haa-v03-2(1113)

**Ocular Comparison:** pl\_blac\_019\_00012-  
1000\_comparisons.txt

**Ocular DIPL XML:** pl\_blac\_019\_00012-1000\_dipl.alto.xml

**Ocular NORM XML:** pl\_blac\_019\_00012-1000\_norm.alto.xml

**Ocular Normalized Transcription:** pl\_blac\_019\_00012-  
1000\_transcription\_normalized.txt

**Ocular Transcription:** pl\_blac\_019\_00012-  
1000\_transcription.txt

Page Num: 13

ID: 101891 Image: pl\_blac\_019\_00013-1000.jpg

fb\_page\_id: pl\_blac\_019\_00013

pl\_blac\_019-ocr-7.19.17-  
mvf(1325)

**Ocular Comparison:** pl\_blac\_019\_00013-  
1000\_comparisons.txt

**Ocular DIPL XML:** pl\_blac\_019\_00013-1000\_dipl.alto.xml

**Ocular NORM XML:** pl\_blac\_019\_00013-1000\_norm.alto.xml

**Ocular Normalized Transcription:** pl\_blac\_019\_00013-  
1000\_transcription\_normalized.txt

**Ocular Transcription:** pl\_blac\_019\_00013-  
1000 transcription.txt

# In conclusion...

---

- In order to change the perception of the state of the OCR problem, humanists must be aware of and emphasize the poor state of machine-reading for anything other than documents created using modern printing practices.
- UNICODE.
- Tools like the DH Dashboard need to de-privilege specific OCR workflows by allowing for flexibility in terms of both processes and data.

# Thank you!

---

Bryan Tarpley

Initiative for Digital Humanities, Media, and Culture

Texas A&M University

[bptarpley@tamu.edu](mailto:bptarpley@tamu.edu)