

# USING DEEP LEARNING TECHNIQUES TO IMPROVE OPTICAL CHARACTER RECOGNITION FOR HISTORICAL DOCUMENTS

by Mengyue Zhao

[beczhaozmy@gmail.com](mailto:beczhaozmy@gmail.com)

# Problem Introduction

Historians have a hard time extracting large amounts of data from raw scans. Especially from historical tables which commercial for which commercial software are not designed.

Linking different historical sources requires very high accuracy of picture-to-data-frame recognition.

I found a very interesting large dataset that requires an improvement on such software. This project is designed to handle tens of thousands of raw scans as shown in the picture on the right side.

The alternative to this solution is manual input by typists, which is costly, inaccurate and often-time yields unknown or stringent mistakes that are hard to correct.

Generated on 2018-06-30 11:42 GMT / http://hdl.handle.net/2027/loc.ark:/13960/rd47hm1m  
Public Domain / http://www.hathitrust.org/access\_use#pd

8* POST OFFICE DEPARTMENT.				
POST OFFICES—Alabama.				
Post office.	County.	Postmaster.	Compensation.	Balance due U. S.
Blakely.....	Baldwin.....	J. H. Stammyres.....	.....	.....
Blake's Ferry.....	Randolph.....	J. J. Bradley.....	.....	\$17 78
Bloeker's.....	Tuscaloosa.....	William C. Blocker.....	.....	.....
Blount Spring.....	Blount.....	B. H. Sapp.....	.....	.....
Blountsville.....	Blount.....	John L. Hopkins.....	.....	66 38
.....	.....	H. A. Gillespie.....	.....	33 81
Blowing Spring.....	Limestone.....	J. Martin.....	.....	.....
Blue Lick.....	Franklin.....	A. T. Bonds.....	.....	.....
Blue Mountain.....	Calhoun.....	J. W. Dundy.....	.....	.....
Blue Pond.....	Cherokee.....	A. Poore.....	20 29	.....
Blue Spring.....	Morgan.....	C. Ratliff.....	65	.....
Bluffport.....	Sumter.....	W. H. Sanders.....	.....	.....
Bluff Spring.....	Talladega.....	M. Street.....	50 79	.....
Boligee.....	Greene.....	S. L. Gwyn.....	8 85	.....
Bowden.....	Talladega.....	J. W. Barton.....	.....	.....
Bradford.....	Coosa.....	Thomas Harris.....	25 36	.....
Braggs.....	Lowndes.....	R. D. Twombly.....	.....	.....
Branchville.....	St. Clair.....	S. G. Hurst.....	32 26	.....
Breckinridge.....	Conecuh.....	W. W. Johnson.....	.....	.....
Brewersville.....	Sumter.....	S. J. Arrington.....	.....	.....
Brickville.....	Lawrence.....	J. N. McCarley.....	20 86	.....
Bridgeport.....	Jackson.....	William C. Glover.....	19 21	.....
Bridge town.....	Shelby.....	Thomas L. Pledger.....	14 23	.....
.....	.....	M. L. Inzer.....	3 27	.....
Bridgeville.....	Pickens.....	M. W. Stephenson.....	23 66	.....
Broken Arrow.....	St. Clair.....	F. B. Walker.....	.....	.....
Brooklyn.....	Conecuh.....	Y. S. Hirsfelder.....	26 65	.....
Brooksville.....	Blount.....	H. C. Reed.....	14 64	.....
Broomtown.....	Cherokee.....	J. Mills.....	.....	.....
Brownsville.....	Talladega.....	Daniel Brown.....	3 38	.....
Brundidge.....	Pike.....	Henry Lane.....	45 03	.....
Brush Creek.....	Perry.....	G. C. Collier.....	13 07	.....
Buchanan.....	Randolph.....	J. Massengale.....	.....	.....
Buckhorn.....	Pike.....	A. R. Fincher.....	.....	.....
Buena Vista.....	Monroe.....	J. P. D. Kelley.....	34 27	.....
Bufoad.....	Barbour.....	M. Patterson.....	.....	.....
Bulger's Mills.....	Coosa.....	J. H. Hubbard.....	11 12	.....
Bullock.....	Coffee.....	J. Thomas.....	.....	.....
Burleson.....	Franklin.....	A. H. Justice.....	.....	.....
Burnsville.....	Franklin.....	A. A. Hughes.....	.....	.....
Burnt Corn.....	Dallas.....	Joshua Dunn.....	.....	.....
Bushville.....	Monroe.....	J. N. Dennard.....	.....	.....
Butler.....	Barbour.....	T. C. Parker.....	119 53	.....
.....	Choctaw.....	J. T. Foster.....	5 81	.....
Butler Spring.....	Butler.....	Chas. C. Hill.....	.....	.....
Buyckville.....	Butler.....	T. C. Watts.....	.....	.....
Buzbeville.....	Coosa.....	W. H. Spigner.....	2 31	.....
Cababa.....	Coffee.....	Wm. F. Johnson.....	354 32	.....
Cainland.....	Dallas.....	R. J. Travers.....	19 85	.....
Calhoun.....	Benton.....	Richard Jenkins.....	43 10	.....
Cambridge.....	Lowndes.....	A. McCaskill.....	.....	.....
Camden.....	Dallas.....	C. M. Cochran.....	411 07	.....
Campbell's Home.....	Wilcox.....	T. C. Brewer.....	12 15	.....
Camp Hill.....	Shelby.....	H. B. Jennings.....	19 47	.....
Camp Spring.....	Blount.....	J. H. Campbell.....	.....	.....
Carlowesville.....	Tallapoosa.....	B. Conire.....	.....	.....
Carrollton.....	Camp Spring.....	William A. Wilam.....	.....	.....
Carthage.....	Dallas.....	J. D. Alison.....	.....	.....
Catoma.....	Carrollton.....	William C. Wilson.....	153 21	.....
Cave Spring.....	Carthage.....	T. Brown.....	13 69	.....
Cedar Bluff.....	Montgomery.....	Benjamin Rodgers.....	10 63	.....
Cedar Grove.....	Fayette.....	S. Hamill.....	1 97	.....
.....	Cherokee.....	John Lawrence.....	9 34	.....
Cedar Plains.....	Jefferson.....	J. N. Benkett.....	8 84	.....
Cedar Springs.....	Morgan.....	N. W. McDaniel.....	.....	.....
Central Institute.....	Benton.....	J. Herring.....	15 06	.....
Centre.....	Coosa.....	William Bell.....	18 20	.....
Centre Hill.....	Limestone.....	R. T. Smith.....	.....	.....
Centre Star.....	Lauderdale.....	W. H. Stiff.....	4 44	.....
Centerville.....	Lawrence.....	William R. Crutcher.....	6 56	.....
Clark Bluff.....	Bibb.....	B. P. Joiner.....	10 53	.....
Chambers C. H.....	Marion.....	J. S. Kennedy.....	19 6	.....
.....	Chambers.....	E. B. Ford.....	88 17	.....
Chanahatchee.....	.....	John A. Frazer.....	96 22	.....
Chandler's Springs.....	Tallapoosa.....	J. R. Brooks.....	26 09	.....
Chepittee.....	Talladega.....	S. D. Watson.....	21 24	.....
Chesnut Creek.....	Blount.....	W. B. Deavor.....	.....	.....
.....	Autauga.....	A. J. Cooper.....	.....	.....

Digitized by  
INTERNET ARCHIVE

Original from  
LIBRARY OF CONGRESS

Tesseract is an optical character recognition software initially maintained by HP in the 1990s, later open-sourced and now maintained by Google.

**What is Tesseract and how does it work?**

## POST OFFICES—Alabama.

Post office.	County.	Postmaster.	Compensation.	Balance due U. S.
Thorn Hill .....	Walker .....	Orrin Davis .....		
Toll Gate .....	Marion .....	J. D. White .....		
Toloca .....	Butler .....	William M. Knight .....		\$22 90
Tompkinsville .....	Choctaw .....	William H. Simmons .....		2 78
Town Creek .....	Lawrence .....	T. B. Hazlewood .....		
Traveller's Rest .....	Coosa .....	S. A. Pearce .....		30 91
Trenton .....	Jackson .....	E. C. Williams .....		
Triana .....	Madison .....	William B. Smith .....		17 53
Trion .....	Tuscaloosa .....	J. W. White .....		
Trout Creek .....	St. Clair .....	H. W. Box .....		
Troy .....	Pike .....	B. T. Murphree .....		79 84
Truss .....	Jefferson .....	T. K. Truss .....		
Turkeytown .....	Cherokee .....	R. S. Cowan .....		65 59
Tusculum .....	Franklin .....	Samuel Finley .....		225 86
Tusculum .....	Tuscaloosa .....	William L. Maunz .....		117 87
Tuskegee .....	Macon .....	P. G. Guild .....		315 55
Uchee .....	Russell .....	John Howard .....		570 22
Union .....	Greene .....	R. E. Covington .....		37 30
Union Springs .....	Macon .....	William Miller .....		116 91
Uniontown .....	Perry .....	A. McGlaty .....		454 73
Upham .....	Macon .....	P. A. Booker .....		236 38
Valherose Springs .....	Morgan .....	R. C. Benson .....		
Valley Head .....	De Kalb .....	J. J. Giers .....		16 56
Van Buren .....	De Kalb .....	William O. Winston .....		19 45
Vienna .....	Pickens .....	M. C. Newman .....		2 98
Village Springs .....	Blount .....	H. L. C. Connerly .....		
Vilula .....	Russell .....	A. P. Gillespie .....		13 89
Vineville .....	Jefferson .....	J. C. Norman .....		38 89
Viole .....	Blount .....	J. Vines, Jr. .....		12 03
Waco .....	Franklin .....	Joel Blackburn .....		
Wadrop's Mill .....	Jefferson .....	R. H. Ransom .....		89 8
Walnut Grove .....	Blount .....	A. G. Waldrop .....		4 23
Walnut Hill .....	Tallapoosa .....	H. Cornelius .....		4 90
Warrenton .....	Marshall .....	William Johnson .....		1 10
Warrior Stand .....	Macon .....	D. Smith .....		86 25
Warsaw .....	Sumter .....	Henry Lane .....		2 07
Water Cure .....	Cherokee .....	Wiley S. Lane .....		
Waterloo .....	Lauderdale .....	Benjamin White .....		
Waverly .....	Chambers .....	S. A. Ross .....		3 46
Weedsway .....	Randolph .....	T. H. Reynolds .....		67 69
Weedsway .....	Randolph .....	D. A. Patrick .....		5 93
Weedsway .....	Randolph .....	R. Lattimer .....		334 33
Weedsway .....	Randolph .....	J. Eichelberger .....		
Weedsway .....	Talladega .....	George Riser .....		
Weedsway .....	Benton .....	P. W. Vaughan .....		63
Wellborn, (c. h.) .....	Coosa .....	J. B. Peacock .....		
Weogufka .....	Pickens .....	John Clift .....		7 56
Wesley Chapel .....	Randolph .....	S. N. Sneed .....		12 39
Westmorelandville .....	Lauderdale .....	J. E. Morefield .....		
Westover .....	Covington .....	R. W. Hawkins .....		24 91
Westville .....	Dale .....	John Hallford .....		10 68
Whistler .....	Mobile .....	M. White .....		33 36
White Plains .....	Pike .....	L. Bryan .....		120 96
White Water .....	Pike .....	E. Fitzgerald .....		86 42
Whitesburg .....	Madison .....	George Borden .....		24 08
Wilcox Springs .....	Wilcox .....	W. J. A. J. Hilliard .....		5 69
Wildwood .....	Randolph .....	William Hanner .....		4 98
Williams's Cross Roads .....	Choctaw .....	J. A. Brown .....		
Williams's Mill .....	Covington .....	J. B. Wood .....		1 97
Will's Valley .....	Cherokee .....	D. A. Perryman .....		
Wilsonville .....	Shelby .....	A. S. Sade .....		5 68
Wied Creek .....	Talladega .....	J. B. Williams .....		11 47
Waterborough .....	Russell .....	R. B. Lankford .....		3 06
Wococochee .....	Russell .....	R. R. King .....		12 71
Wolf Creek .....	St. Clair .....	J. Johnson, Jr. .....		
Woodstock .....	Tuscaloosa .....	L. Brock .....		9 55
Woodville .....	Jackson .....	A. Hammond .....		
Worth .....	De Kalb .....	Elias J. Brock .....		1 05
Yellow Creek .....	Payette .....	James Ray .....		
Yongesborough .....	Russell .....	J. J. Dillard .....		47 03
York .....	Walker .....	H. K. Wile .....		25 90
Youngville .....	Tallapoosa .....	J. B. Morton .....		16 00
Zachary .....	Marshall .....	W. C. Yonge .....		13 69
		Walker .....		
		E. Clark .....		33 83
		R. G. Young .....		
		William W. McCutchen .....		13 29

3\*

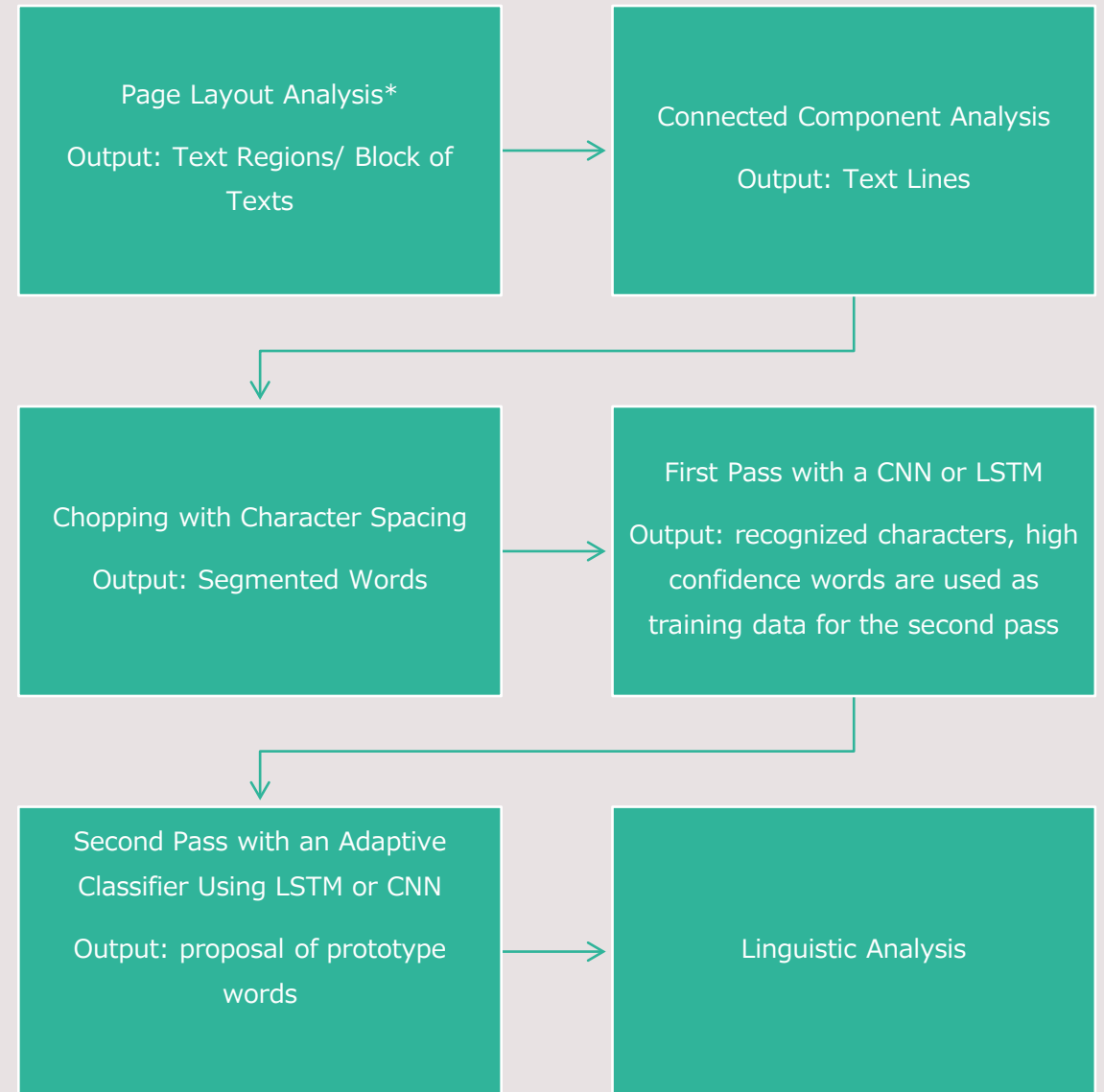
# How does Tesseract (or OCR in general) Work?

Pre-  
processingPage Layout  
AnalysisText Line  
AnalysisCharacter  
SegmentationOptical  
Character  
Recognition  
(OCR)Post  
Processing



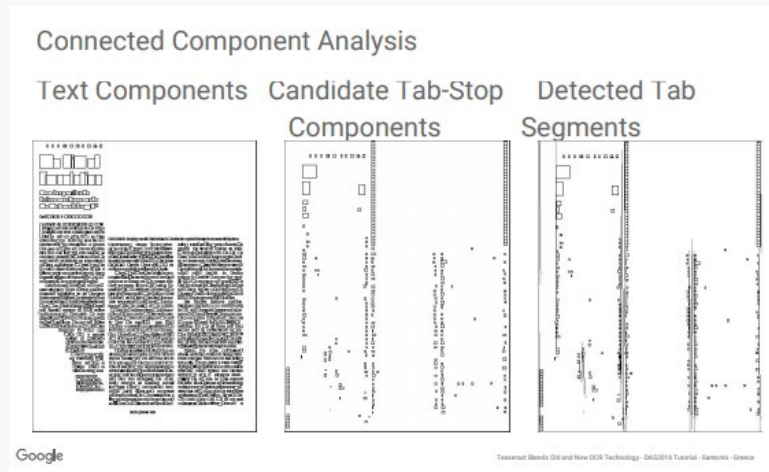
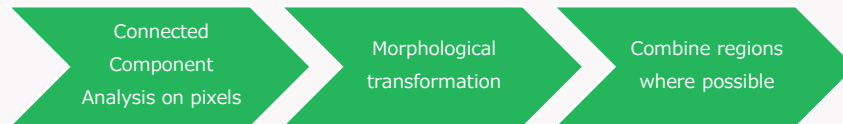


# Tesseract Architecture

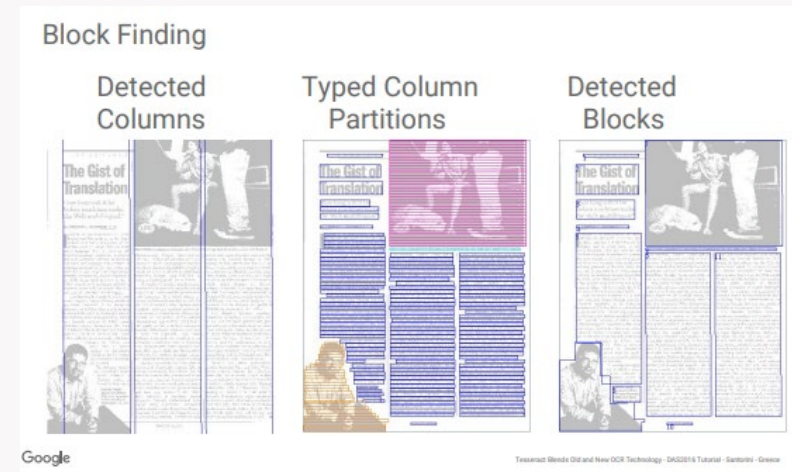
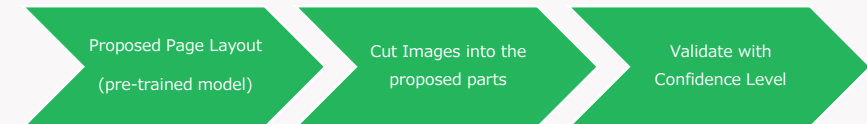


# Page Layout Analysis

## Bottom-up Approach



## Top-down Approach



# Using LSTM in OCR

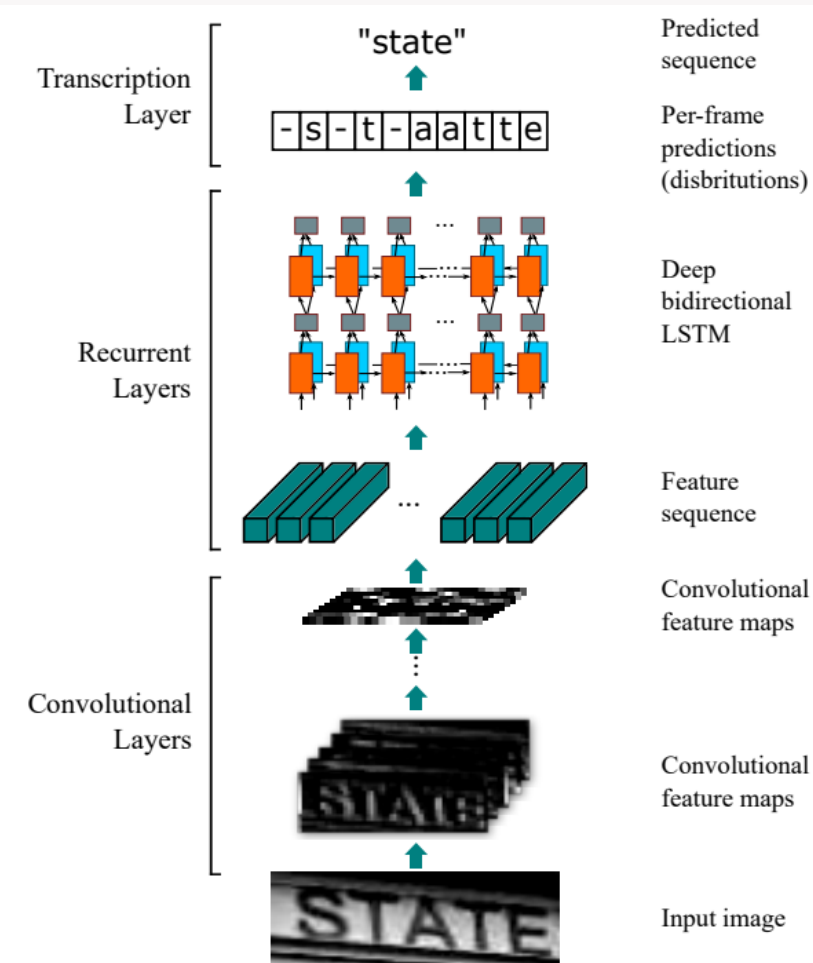


Figure 1. The network architecture. The architecture consists of three parts: 1) convolutional layers, which extract a feature sequence from the input image; 2) recurrent layers, which predict a label distribution for each frame; 3) transcription layer, which translates the per-frame predictions into the final label sequence.

# Links

## Visualized Tesseract Architecture

<https://limitlessdatascience.wordpress.com/2019/07/01/tesseract-4-0-intro-installation/>

## How does Tesseract use LSTM

[https://github.com/tesseract-ocr/docs/blob/master/das\\_tutorial2016/6ModernizationEfforts.pdf](https://github.com/tesseract-ocr/docs/blob/master/das_tutorial2016/6ModernizationEfforts.pdf)

<https://nanonets.com/blog/ocr-with-tesseract/>

## How LSTM work for OCR

<https://towardsdatascience.com/a-gentle-introduction-to-ocr-ee1469a201aa>



# Page Segmentation Solutions

- Left: Tesseract Segmentation Results
- Right: Customized Page Segmentation

[illegible]

PORT OFFICE DEPARTMENT.					
PORT OFFICE - TON.					
No.	Name.	Rank.	Pay.	Allowances.	Total.
1					
2					
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14					
15					
16					
17					
18					
19					
20					
21					
22					
23					
24					
25					
26					
27					
28					
29					
30					
31					
32					
33					
34					
35					
36					
37					
38					
39					
40					
41					
42					
43					
44					
45					
46					
47					
48					
49					
50					
51					
52					
53					
54					
55					
56					
57					
58					
59					
60					
61					
62					
63					
64					
65					
66					
67					
68					
69					
70					
71					
72					
73					
74					
75					
76					
77					
78					
79					
80					
81					
82					
83					
84					
85					
86					
87					
88					
89					
90					
91					
92					
93					
94					
95					
96					
97					
98					
99					
100					

# OCR Results

## (with Customized Segmentation)

Delta	Randolph	T A Cantrell	07
Democrat	Walker	B Rice	
Yemocrat	Marengo	Thomas G W Shannon	634
Jemopolis	Tallapoosa	H P Smith	
rr?ve Soto	Choctaw	T Poole	
Ye Sotoville	Marion	Win W Davidson	184
YYetroit	Macon	AG Tuttle	3209
Dicks Creek	Alacon	Thomas Anthony	680
Dickson	Franklin	T A Beverly	53
lickson	Butler	Jacob Gross	
Delask	Marengo	S Lee	508
YYidask	Marensro	W oO Brannon	
Dixon's Mill	Jackson	J M Warren	16692
Dixons Mill	Dale		
Dodsonville	Russell		
Jodsonville	Lawrence		
Ootham			
Pothan			
Doyer			
Dry Creek			
Dry Creek			



# Proposed Solutions

- Customized Preprocessing [done]
- Page Layout Analysis with Morphological Transformation [done]
- Page Segmentation with Morphological Transformation [done]
- Fixed-width Character Segmentation [tuning]
- Retrain Neural Network on Specific Fonts [done]
- Restrained Classifier with a Customized Dictionary (word list) [doing]
- Restricted Output Base [done]
- Post Processing [doing]