

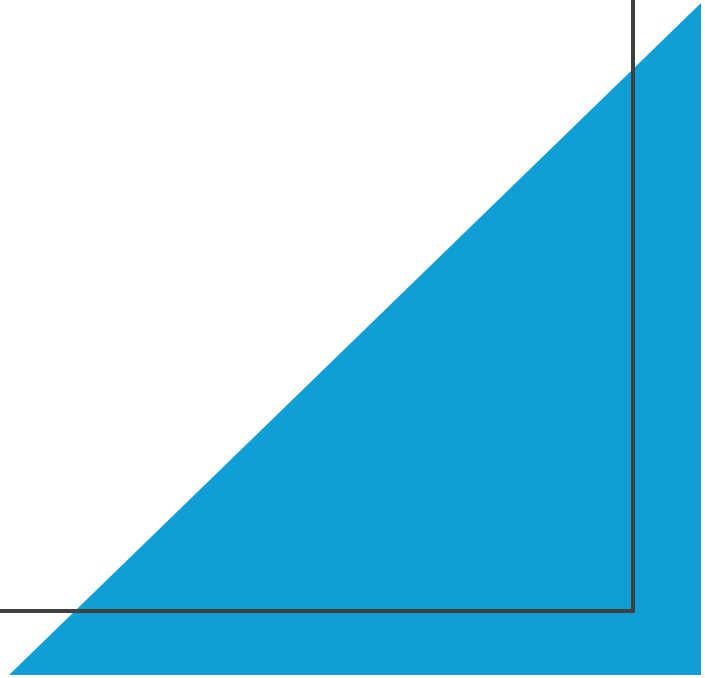
Enhancing the Precision of Vast Historical Datasets via Deep Learning and AI: A Case Study on Mitigating Misinterpretation Challenges within the 1950 U.S. Census

Mengyue Zhao (University of Oxford)

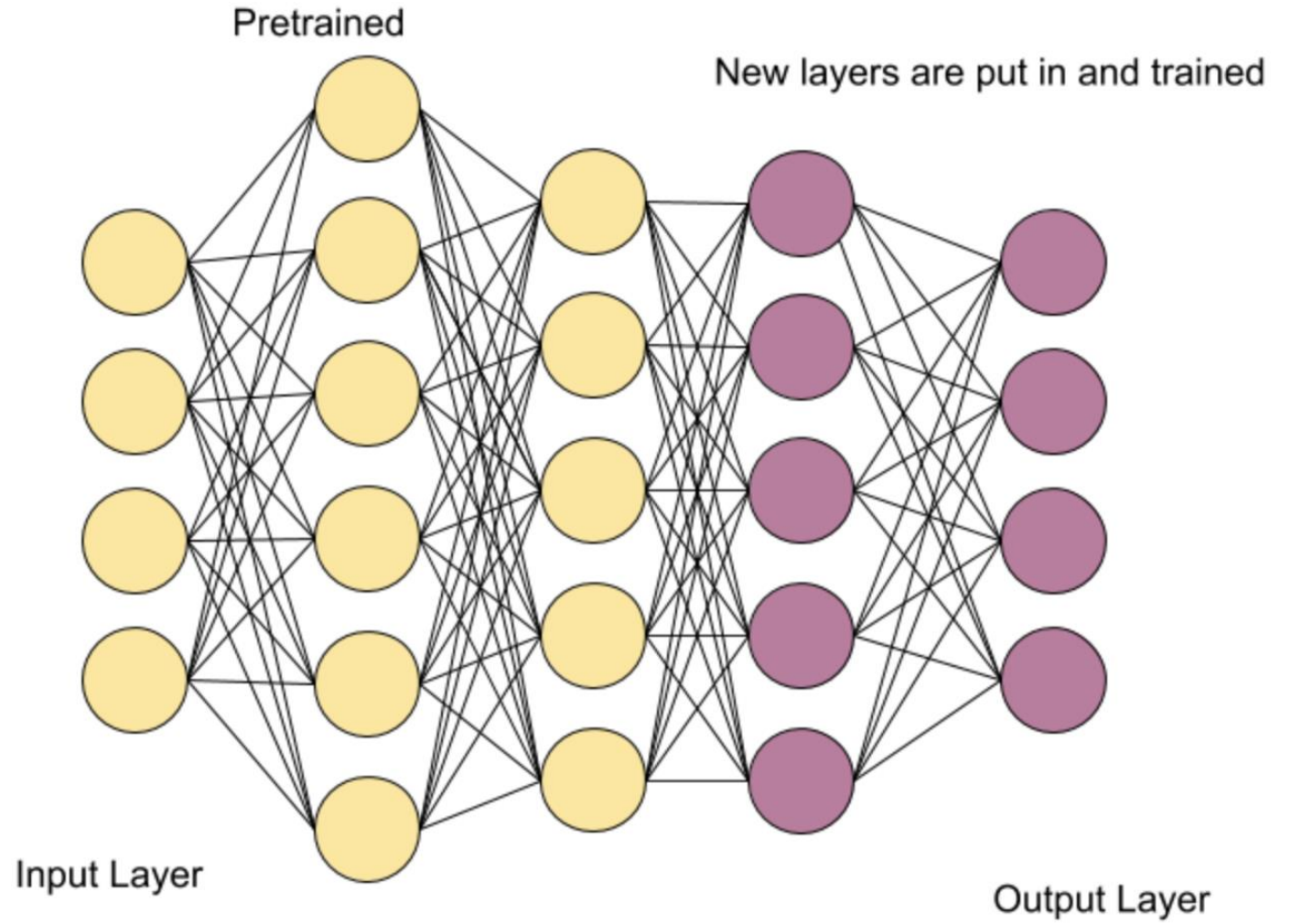
beczhaozmy0@gmail.com

Agenda

- What is deep-learning & transfer learning?
- Deep-Learning Frameworks
 - Machine Learning Iteration Process
 - How robust is this approach?
- Use Cases:
 - Transcribing the Official Register of the United States
 - Reduce 1950 Census Overcounts
- The Importance Page Segmentation
- The role of Generative-AI in Fine-tuning tasks
 - Can it replace traditional ML?



Transfer Learning



The 1950 census data

- an **overcount** of approximately **two million individuals**.
- Source of Error: OCR mistakes
 - “No One” -> “Owen” ;
 - “Not at Home” -> “Norton”
- “Strike-through” Indicators
- Deep-Learning Solution
 - Page Layout Analysis
 - Image Classification



1950 United States Federal Census for Scheel George

California > San Mateo > Burlingame > 41-20

Save

CONFIDENTIAL
 U. S. DEPARTMENT OF COMMERCE
 BUREAU OF THE CENSUS
FORM P1
1950 CENSUS OF POPULATION AND HOUSING

COUNTY: **SAN MATEO**
 INCORPORATED PLACE OR TOWNSHIP: **BURLINGAME CITY**
 U. S. D. NUMBER: **41-20**

ENUMERATOR'S SIGNATURE: *Margie J. Callow*
 NUMBER: **25**

Notes: Not known for sure.
 See SP3-4 3

| FOR HEAD OF HOUSEHOLD | | | | | | FOR ALL PERSONS | | | | | | | | | | FOR PERSONS 14 YEARS OF AGE AND OVER | | | | | | | |
|-----------------------|--------------------------------------|--------------------------------|-----------------------------------|--|---------------------------------|------------------------------------|--------------|------|-----|--|-----------------|---------------|----------------|-------------|------------|--------------------------------------|-----------------|-------------|----|-----------|--------------|-----------|---|
| LINE NUMBER | NAME OF HEAD, SURNAME, OR FIRST NAME | SERIAL NUMBER OF DWELLING UNIT | IS THIS HOUSE ON A FARM OR RANCH? | IS THIS HOUSE ON A PLACE OF THREE OR MORE ACRES? | APPROXIMATE QUANTITATIVE NUMBER | NAME | RELATIONSHIP | RACE | SEX | AGE | MARRIAGE STATUS | BIRTHPLACE | NATURALIZATION | WORK STATUS | OCCUPATION | INDUSTRY | CLASS OF WORKER | LINE NUMBER | | | | | |
| | | | | | | | | | | | | | | | | | | | 1 | 2 | 3 | 4 | 5 |
| | | | | | | LORRAINE K. daughter | W | F | 1 | nr | Calif. | | | | | | | 1 | | | | | |
| | | | | | | GEORGE, CLORINDA | | | | | | | | | | | | 2 | | | | | |
| 1473 | 290 | | not at home see sheet | | | Surname: George Given Name: Scheel | | | | | | | | | | | | | 48 | housework | Private home | P 7208261 | 3 |
| 1467 | 291 | | not at home see sheet | | | | | | | | | | | | | | | 4 | | | | | |
| 1465 | 292 | no | no | | | DODD, JOHN M. | head | W | M | 53 | mar | New York | 02/1 | WK | | | | 5 | | | | | |
| | | | | | | FLORENCE A. | wife | W | F | 50 | mar | New York | 02/1 | H | no | no | no | 6 | | | | | |
| | | | | | | PATRICIA | daughter | W | F | 11 | nr | New York | 02/1 | | | | | 7 | | | | | |
| 1455 | 293 | | not at home see sheet | | 77 lines | | | | | 12, 13, 14, 15, 16, 17, 18, 19, 20, 21 | | | | | | | 8 | | | | | | |
| 1458 | 294 | | not at home see sheet | | 77 lines | | | | | 11 | | | | | | | 9 | | | | | | |
| 1449 | 295 | no | no | | | BIONDI, DINO H. | head | W | M | 40 | mar | Pennsylvania | 02/3 | OT | no | yes | | 10 | | | | | |
| | | | | | | IRENE P. | wife | W | F | 36 | nr | Calif. | | H | no | no | no | 11 | | | | | |
| | | | | | | J. CHRISTIAN | son | W | M | 7 | nr | Calif. | | | | | | 12 | | | | | |
| 1445 | 296 | no | no | | | EMPEY, WALTER T. | head | W | M | 39 | mar | Calif. | | OT | no | no | yes | 13 | | | | | |
| | | | | | | ERMA F. | wife | W | F | 33 | mar | Calif. | | H | no | no | no | 14 | | | | | |
| | | | | | | DARRELL R. | son | W | M | 1 | nr | Calif. | | | | | | 15 | | | | | |
| 1441 | 297 | no | no | | | MARSH, DAVID G. | head | W | M | 68 | mar | West Virginia | 05/5 | OT | no | no | no | 16 | | | | | |
| | | | | | | ANNA C. | wife | W | F | 62 | mar | Maryland | 05/2 | H | no | no | no | 17 | | | | | |
| | | | | | | ROBERT D. | son | W | M | 34 | nr | Pennsylvania | 02/3 | WK | | | | 18 | | | | | |
| 1437 | 298 | no | no | | | VACANT | | | | | | | | | | | 19 | | | | | | |
| 1433 | 299 | no | no | | | GALITZ, EVELYN M. | head | W | F | 41 | mar | Illinois | 03/3 | H | no | no | no | 20 | | | | | |
| | | | | | | RICHARD L. | son | W | M | 13 | nr | Illinois | 03/3 | | | | | 21 | | | | | |
| | | | | | | GORDON W. | son | W | M | 5 | nr | Calif. | | | | | | 22 | | | | | |
| | | | | | | KEITH G. | son | W | M | 5 | nr | Calif. | | | | | | 23 | | | | | |
| 1429 | 300 | no | no | | | VACANT | | | | | | | | | | | 24 | | | | | | |
| 1425 | 301 | no | no | | | HOWELL DELRUS D. | head | W | M | 57 | nr | Calif. | | OT | no | no | no | 25 | | | | | |

3049 4-17
 Cabrillo Terrace

1950 United States Federal Census for Lome See Howell

California > San Mateo > Burlingame > 41-20

Save ▾

more to discover in this
Dive deeper into the life of
See Howell.

Start

HOUSEHOLD

FOR ALL PERSONS

FOR PERSONS 14 YEARS OF AGE AND OVER

| LINE NUMBER | Name of street, avenue, or road | House (and apartment) number | Serial number of dwelling unit | Is this house on a farm (or ranch)? | Is this house on a place of three or more acres? | Agriculture Questionnaire Number | NAME What is the name of the head of this household? What are the names of all other persons who live here? List in this order: The head His wife Unmarried sons and daughters (in order of age) Married sons and daughters and their families Other relatives Other persons, such as lodgers, roomers, maids or hired hands who live in, and their relatives | RELATIONSHIP Enter relationship of person to head of the household, as Head Wife Daughter Grandson Mother-in-law Lodger Lodger's wife Maid Hired hand Patient, etc. | RACE White (W) Negro (Neg) American Indian (Ind) Japanese (Jap) Chinese (Chi) Filipino (Fil) Other race—spell out | SEX Male (M) Female (F) | How old was he on his last birthday? (If under one year of age, enter month of birth as April, May, Dec., etc.) | Is he now married, widowed, divorced, separated, or never married? (Mar. W, D, Sep. N, or V) | What State (or foreign country) was he born in? If born outside Continental United States, enter name of Territory, possession, or foreign country Distinguish Canada-French from Canada-other | If foreign born— Is he naturalized? (Yes, No, or AP for born abroad of American parents) | What was this person doing most of last week—working, keeping house, or something else? (Wk., H, O, or U for unable to work) | If H or O in item 15—Did this person do any work at all last week, not counting work around the house? (Include work for pay, in own business, profession, on farm, or unpaid family work) (Yes or No) | If No in item 16—Was this person looking for work? (See Special Cases below) (Yes or No) | If No in item 17—Even though he didn't work last week, does he have a job or business? (Yes or No) | If Wk in item 15 or Yes in item 16—How many hours did he work last week? (Include unpaid work on farm or business) (Number of hours) | 1. If employed (Wk in item 15, or Yes in item 16) 2. If looking for work (Yes in item 17), describe in 3. For all other persons, leave blank What kind of work was he doing? For example: Nails heels on shoes..... Chemistry professor..... Farmer..... Farm helper..... Armed forces..... Never worked..... (Occupation) | |
|-------------|---------------------------------|------------------------------|--------------------------------|-------------------------------------|--|----------------------------------|--|---|--|-------------------------------|--|---|--|--|---|--|--|--|--|--|------------------------------|
| | | | | | | | | | | | | | | | | | | | | | LEAVE BLANK |
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20a | | |
| | | | | | | LESLEE ANN | daughter | W | F | 11 | never | Calif. | | | | | | | | | |
| | | 1467 | 303 | | | not at home see | | | | | | | | | | | | | | | |
| | | 1405 | 304 | no | no | BEGGS, WILLIAM H. | head | W | M | 63 | mar | Northern Ireland | 11 | yes | WK | | | | | 32 | Carpenter |
| 4 | | | | | | ELIZABETH L. | wife | W | F | 65 | mar | Northern Ireland | 13 | yes | H | no | no | no | | | |
| | | 1800 | 305 | no | no | CUTSHAW, ROBERT P. | head | W | M | 35 | mar | North Carolina | 05 | | WK | | | | | 65 | Contractor division salesman |
| | | | | | | LUCY M. | wife | W | F | 34 | mar | Tennessee | 06 | 2 | H | no | no | no | | | |
| | | | | | | ROBERT ANDREW | son | W | M | 4 | never | Tennessee | 06 | 2 | # | | | | | | |
| | | 1804 | 306 | no | no | RUSCH, CLARENCE E. | head | W | M | 40 | mar | Wisconsin | 35 | | WK | | | | | 60 | Proprietor |
| | | | | | | ELINOR M. | wife | W | F | 32 | mar | Calif. | | | H | no | no | yes | | | bookkeeping |
| | | | | | | JUDITH A. | daughter | W | F | 9 | never | Calif. | | | | | | | | | |
| | | | | | | SELMA S. | daughter | W | F | 7 | never | Calif. | | | | | | | | | |
| | | | | | | MARTHA E. | daughter | W | F | 4 | never | Calif. | | | | | | | | | |
| | | | | | | EMIL | father | W | M | 64 | wa | Wisconsin | 08 | 5 | WK | | | | | 60 | repair man |
| | | 1810 | 307 | no | no | not at home see sheet | | | | | | | | | | | | | | | |
| | | 1406 | 308 | no | no | not at home see sheet | | | | | | | | | | | | | | | |

Surname: Howell Given Name: Lome See

6 sep 4-12

Willie

27 of 45 lines 12, 13, 14, 15, 16

lines 17, 18, 19, 20

1950 United States Federal Census for Thomas Monahan

California > Los Angeles > Los Angeles > 66-1560

Save

| Line Num | Street Ne | House of | Dwelling | Farm | Acres | Question | Name | Relations | Code A | Race | Gender | Age | Marit | Birth Place | Code B | Citizensh | Occu | Work | Seeking | Employ | Hour | Occupation | Industry | Worker C | Code C | Code C | Code C | LINE NUMBER |
|----------|-----------|----------|----------|------|-------|----------|-----------------------|-----------|--------|------|--------|-----|-------|---------------|--------|-----------|------|------|---------|--------|------|-------------------------|-------------------------|----------------------|--------|--------|--------|-------------|
| 15 | | | No | No | | | Linn, Lewis F. | Head | | W | M | 66 | Mar | New York | 021 | | WK | | | | 28 | music teacher | retail music store | P | 057 | 698 | 1 | |
| 4719 | | | 14 | No | No | | - Beatrice (now) | wife | | W | F | 71 | Mar | New York | 021 | | H | Yes | | | 6 | music teacher | dance school | P | 057 | 888 | 2 | |
| 4723 | | | | | | | No one at home | see sheet | | | 73 | | | | | | | | | | | | | | | | 3 | |
| 4735 | | | | | | | No one at home | see sheet | | | 72 | | | | | | | | | | | | | | | | | 4 |
| 4807 | | | 17 | No | No | | No one at home | see sheet | | | 72 | | | | | | | | | | | | | | | | | 5 |
| | | | 18 | No | No | | No one at home | see sheet | | | 71 | | | | | | | | | | | | | | | | | 6 |
| | | | 19 | No | No | | Louisa, Helen P. | Head | | W | F | 54 | WD | Canada-Other | 161 | Yes | WK | | | | 48 | mail operator | hotel | P | 790 | 836 | 7 | |
| | | | | | | | David J. | son | | W | M | 27 | Mar | Massachusetts | 014 | | WK | | | | 40 | operator | Veterans Administration | G | 341 | 916 | 8 | |
| | | | | | | | McCabe, Mary E. | sister | | W | F | 62 | WD | Canada-Other | 161 | Yes | H | No | Yes | No | | | sales clerk | retail ready to wear | P | 490 | 656 | 9 |
| | | | | | | | Monahan, Kathryn A. | sister | | W | F | 64 | WD | Canada-Other | 161 | Yes | WK | | | | 40 | mail operator | hotel | P | 790 | 836 | 10 | |
| 4551 1/2 | | | 20 | No | No | | No one at home | see sheet | | | 71 | | | | | | | | | | | | | | | | | 11 |
| 4551 | | | 21 | No | No | | Vacant | | | | | | | | | | | | | | | | | | | | | 12 |
| 4539 | | | 22 | No | No | | No one at home | | | | | | | | | | | | | | | | | | | | | 13 |
| 4529 | | | 23 | No | No | | Vacant | | | | | | | | | | | | | | | | | | | | | 14 |
| 4418 | | | 24 | No | No | | Tomplins, Clifford M. | Head | | W | M | 52 | Mar | Canada-Other | 161 | Yes | WK | | | | 60 | proprietor | retail grocery store | O | 290 | 636 | 15 | |
| | | | | | | | Jean M. | wife | | W | F | 48 | Mar | Canada-Other | 161 | No | WK | | | | 40 | sales clerk | retail grocery | P | 490 | 636 | 16 | |
| | | | | | | | Wallace M. | son | | W | M | 22 | Mar | California | | | OT | No | No | No | | | | | | | | 17 |
| | | | | | | | Shanon E. | daughter | | W | F | 23 | Mar | Calif. | | | WK | | | | 40 | nurse | county hospital | G | 058 | 869 | 18 | |
| 4464 | | | 25 | No | No | | No one at home | see sheet | | | 71 | | | | | | | | | | | | | | | | | 19 |
| 4472 | | | 26 | No | No | | Acosta, Joe A. | Head | | W | M | 24 | Mar | Calif. | | | OT | No | Yes | | | | spray operator | pottery factory | P | 490 | 319 | 20 |
| | | | | | | | Jean (now) | wife | | W | F | 22 | Mar | Calif. | | | H | No | Yes | | | | counter girl | retail bakery | P | 490 | 636 | 21 |
| | | | | | | | Frank A. | son | | W | M | 5 | Mar | Calif. | | | | | | | | | | | | | | 22 |
| 4472 | | | | | | | Sanchez, Pete L. | lodger | | W | M | 43 | D | Mexico | 262 | No | OT | No | Yes | | | | conductor | urban transportation | P | 631 | 516 | 23 |
| 4478 | | | 27 | No | No | | Walker, Raymond E. | Head | | W | M | 29 | Mar | Missouri | 043 | | WK | | | | 54 | assistant store manager | retail sales | P | 090 | 699 | 24 | |
| | | | | | | | Helen J. | wife | | W | F | 27 | Mar | Missouri | 043 | | H | No | No | No | | | | | | | | 25 |
| | | | | | | | Paula J. | daughter | | W | F | 4 | Mar | Missouri | 043 | | | | | | | | | | | | | 26 |
| 4478 1/2 | | | 28 | No | No | | Hartman, Emory J. | Head | | W | M | 34 | Mar | Illinois | 033 | | OT | No | Yes | | | | bar tender | bar | P | 750 | 679 | 27 |
| | | | | | | | Francis M. | wife | | W | F | 26 | Mar | Arkansas | 071 | | H | No | No | No | | | | | | | | 28 |
| | | | | | | | Shanna M. | daughter | | W | F | 7 | Mar | Calif. | | | | | | | | | | | | | | 29 |
| | | | | | | | Gary F. | son | | W | M | 5 | Mar | Calif. | | | | | | | | | | | | | | 30 |
| | | | | | | | Kristin L. | daughter | | W | F | 3 | Mar | Calif. | | | | | | | | | | | | | | 31 |

Surname: Monahan Given Name: Thomas

HOUSEHOLD CONTINUED ON NEXT SHEET

Notes

Size of the Prize

- overcounts across all Census Collections

U.S. DEPARTMENT OF COMMERCE
1950 CENSUS OF POPULATION AND HOUSING

Arizona
Santa Cruz
Nozales City
12-3

APRIL 6
C. C. [unclear]

CONFIDENTIAL

FOR ALL PERSONS

| FOR HEAD OF HOUSEHOLD | | | | | | FOR ALL PERSONS | | | | | | | | | | | | | FOR PERSONS 14 YEARS OF AGE AND OVER | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|-----------------------|-----------------------------------|-----|-----|---------------|----------------|-----------------|----------|---------------|----------------|------|------------|---------------|----------------|-----|-----|---------------|----------------|---------------|--------------------------------------|---------------|----------------|------|----------|---------------------|----------------|-----|-----|---------------|----------------|----------|----------|---|----|------|----------|---------------------|------|---|----|------|----------|-----------|------|---|----|------|----------|------------------|------|---|----|------|---------|----------|------|---|----|------|---------|-------------|-----|---|----|------|---------|--------|-----|---|----|------|---------|-----------------|------|---|----|------|--------|------------|------|---|----|------|--------|-------|----------|---|----|------|---------|--------|----------|---|----|------|---------|
| NAME | RELATIONSHIP TO HEAD OF HOUSEHOLD | SEX | AGE | DATE OF BIRTH | ETHNIC OR RACE | SEX | AGE | DATE OF BIRTH | ETHNIC OR RACE | SEX | AGE | DATE OF BIRTH | ETHNIC OR RACE | SEX | AGE | DATE OF BIRTH | ETHNIC OR RACE | SEX | AGE | DATE OF BIRTH | ETHNIC OR RACE | SEX | AGE | DATE OF BIRTH | ETHNIC OR RACE | SEX | AGE | DATE OF BIRTH | ETHNIC OR RACE | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| LWEIG, FREDERICK | SON | M | 18 | 1932 | ARIZONA | LWEIG, ROBERT | DAUGHTER | F | 17 | 1933 | COSTA RICA | LARA, EUTQUIA | MAID | F | 38 | 1912 | MEXICO | BARNET, FANNY | HEAD | F | 77 | 1853 | MISSOURI | FRAZIER, CHARLES T. | SON-IN-LAW | M | 60 | 1888 | MISSOURI | EMMA MAE | DAUGHTER | F | 49 | 1908 | MISSOURI | HATTENBACH, IRON M. | HEAD | M | 38 | 1912 | COLORADO | ELIZABETH | WIFE | F | 38 | 1912 | COLORADO | JAABERRA, JOSEPH | HEAD | M | 42 | 1908 | ARIZONA | MARGARET | WIFE | F | 36 | 1914 | ARIZONA | JR., JOSEPH | SON | M | 17 | 1933 | ARIZONA | ROBERT | SON | M | 15 | 1935 | ARIZONA | GARCIA, LUIS R. | HEAD | M | 43 | 1907 | MEXICO | IGNACIO J. | WIFE | F | 44 | 1906 | MEXICO | NORMA | DAUGHTER | F | 19 | 1931 | ARIZONA | GLORIA | DAUGHTER | F | 18 | 1932 | ARIZONA |

48 HOUSEHOLDS GENERAL HOUSEHOLD PRIVATE HOUSEHOLD P 720

49 AGENT WHOLESALE GROCERIES PRODUCTS P 291

48 SUPERVISOR BANK P 291

48 IMMIGRATION INSPECTOR U.S. IMMIGRATION SERVICE G 216

45 BUS DRIVER INTER STATE TRANSPORTATION P 465

54 PROPRIETOR RETAIL GROCERY STORE O 276

30 SALES PERSON RETAIL GROCERY STORE NP 746

16 SALES PERSON RETAIL GROCERY STORE NP 746

21 HOUSE FOUND TO BE VACANT AT LATER DATE. DWELLING UNIT INFORMATION ON SHEET 74 LINE 14

21 HOUSE FOUND TO BE VACANT AT LATER DATE. DWELLING UNIT INFORMATION ON SHEET 73 LINE 1

THE QUESTIONS BELOW ARE FOR PERSONS LISTED ON SAMPLE LINES

FOR ALL AGES

FOR PERSONS 14 YEARS OF AGE AND OVER

What country and State was he living in a year ago?

What country was his father and mother born in?

What year, month, and day did you last work for any one?

What kind of work was it?

What kind of business or industry was it?

What kind of work did you do in the last job?

What kind of business or industry was it?

What kind of work did you do in the last job?

What kind of business or industry was it?

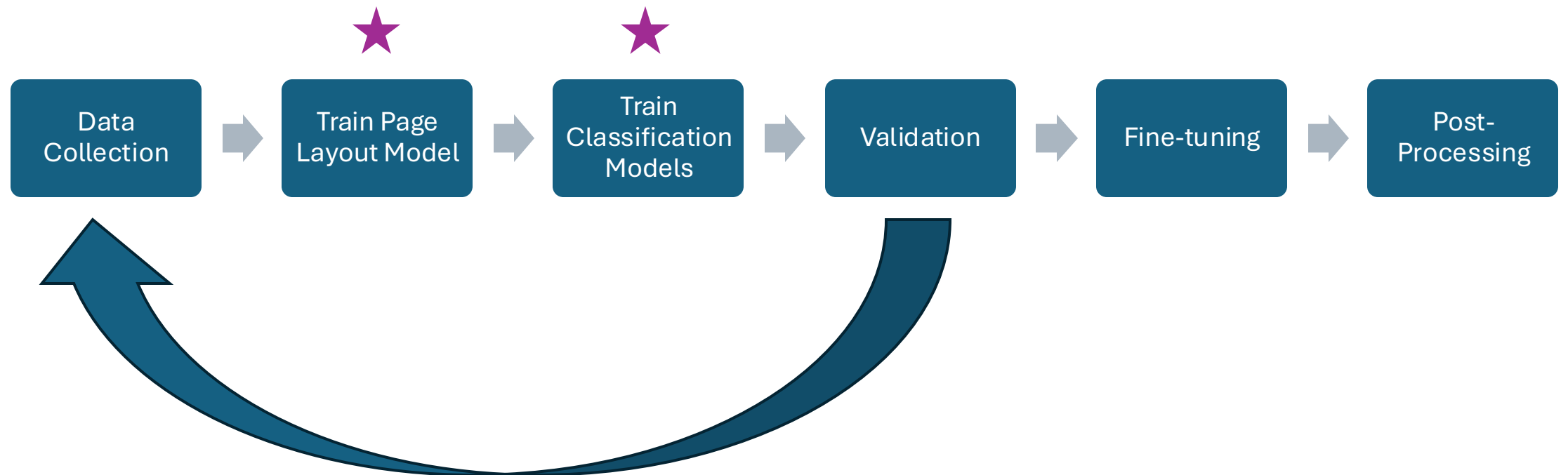
What kind of work did you do in the last job?

What kind of business or industry was it?

What kind of work did you do in the last job?

What kind of business or industry was it?

Machine Learning Iteration Process



| | | | | |
|-----------------------|-----------------|------------------------|-----|------|
| Blowing Spring..... | Limestone..... | H. A. Gillespie..... | 33 | 81 |
| Blue Lick..... | Franklin..... | J. Martin..... | | |
| Blue Mountain..... | Calhoun..... | A. T. Bonds..... | | |
| Blue Pond..... | Cherokee..... | J. W. Dendy..... | | |
| Blue Spring..... | Morgan..... | A. Poore..... | 20 | 20 |
| Bluffport..... | Sumter..... | C. Ratliff..... | | 65 |
| Bluff Spring..... | Talladega..... | W. H. Sanders..... | | |
| Bolige..... | Greene..... | M. Street..... | 20 | 79 |
| Bowden..... | Talladega..... | S. L. Gwyn..... | | 8 85 |
| Bradford..... | Coosa..... | J. W. Barton..... | | |
| Braggs..... | Lowndes..... | Thomas Harris..... | 25 | 36 |
| Branchville..... | St. Clair..... | R. D. Twombly..... | | |
| Breckinridge..... | Conecuh..... | S. G. Hurst..... | 32 | 26 |
| Brewersville..... | Sumter..... | W. W. Johnson..... | | |
| Brickville..... | Lawrence..... | S. J. Arrington..... | | |
| Bridgeport..... | Jackson..... | J. N. McCarley..... | 20 | 86 |
| Bridgetown..... | Shelby..... | William C. Glover..... | 19 | 23 |
| | | Thomus L. Pledger..... | 14 | 23 |
| | | M. L. Inzer..... | | 3 27 |
| Bridgeville..... | Pickens..... | M. W. Stephenson..... | 33 | 06 |
| Broken Arrow..... | St. Clair..... | F. B. Walker..... | | |
| Brooklyn..... | Conecuh..... | Y. S. Hirshfelder..... | 26 | 85 |
| Brooksville..... | Blount..... | H. C. Reed..... | 14 | 64 |
| Broomtown..... | Cherokee..... | J. Mills..... | | |
| Brownsville..... | Talladega..... | Daniel Brown..... | 3 | 38 |
| Bruceville..... | Pike..... | Henry Lane..... | | |
| Brundidge..... | Pike..... | G. C. Collier..... | 45 | 03 |
| Brush Creek..... | Perry..... | J. Massengale..... | 13 | 07 |
| Buchanan..... | Randolph..... | A. B. Fincher..... | | |
| Buckhorn..... | Pike..... | J. P. D. Kelley..... | | |
| Buena Vista..... | Monroe..... | M. Patterson..... | 34 | 27 |
| Buford..... | Barbour..... | J. Thomas..... | | |
| Bulger's Mills..... | Coosa..... | J. H. Hubbard..... | 11 | 12 |
| Bullock..... | Coffee..... | A. H. Justice..... | | |
| Burleson..... | Franklin..... | A. A. Hughes..... | | |
| Burnsville..... | Dallas..... | Josiah Dunn..... | | |
| Burnt Corn..... | Monroe..... | J. N. Dennard..... | | |
| Bushville..... | Barbour..... | T. C. Parker..... | | |
| Butler..... | Choctaw..... | J. T. Foster..... | 119 | 53 |
| | | Chas. C. Hill..... | | 5 81 |
| Butler Spring..... | Butler..... | T. C. Watts..... | | |
| Buycksville..... | Coosa..... | W. H. Spigener..... | | |
| Buzbeeville..... | Coffee..... | Wm. F. Johnson..... | 2 | 31 |
| Cababa..... | Dallas..... | R. J. Travers..... | 354 | 32 |
| Cainland..... | Benton..... | Richard Jenkins..... | 3 | 85 |
| Calhoun..... | Lowndes..... | A. McCaskill..... | 43 | 10 |
| Cambridge..... | Dallas..... | C. M. Cochran..... | | |
| Camden..... | Wilcox..... | T. C. Brewer..... | 411 | 07 |
| Campbell's Home..... | Shelby..... | H. B. Jennings..... | | |
| Campbell's Store..... | Blount..... | J. H. Campbell..... | 12 | 15 |
| Camp Hill..... | Tallapoosa..... | B. Conire..... | 19 | 47 |
| Camp Spring..... | Lawrence..... | William A. Wilam..... | | |
| Carlowsville..... | Dallas..... | J. D. Alison..... | | |
| Carrollton..... | Pickens..... | William C. Wilson..... | 155 | 21 |
| Carthage..... | Tuscaloosa..... | T. Brown..... | 13 | 69 |
| Catoma..... | Montgomery..... | Benjamin Rodgers..... | 10 | 63 |
| Cave Spring..... | Fayette..... | S. Hamil..... | | |
| Cedar Bluff..... | Cherokee..... | John Lawrence..... | 80 | 03 |
| Cedar Grove..... | Jefferson..... | J. N. Benkett..... | 1 | 97 |
| | | N. W. McDaniel..... | 9 | 34 |

Use Case 1: Transcribing *the Official Register of the United States*



Page Segmentation Results

| | | | | |
|----------------|------------|--------------------|-------|---------|
| Blakely | Baldwin | J. H. Stanmyres | | |
| Blake's Ferry | Randolph | J. J. Bradley | | \$17 78 |
| Blocker's | Tuscaloosa | William G. 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. Dendy | | |
| Blue Pond | Cherokee | A. Poore | 20 20 | |
| Blue Spring | Morgan | C. Ratliff | 65 | |
| Bluffport | Sumter | W. H. Sanders | | |
| Bluff Spring | Talladega | M. Street | 90 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 | |

| | | | | |
|----------------|------------|--------------------|-------|---------|
| Blakely | Baldwin | J. H. Stanmyres | | |
| Blake's Ferry | Randolph | J. J. Bradley | | \$17 78 |
| Blocker's | Tuscaloosa | William G. 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. Dendy | | |
| Blue Pond | Cherokee | A. Poore | 20 20 | |
| Blue Spring | Morgan | C. Ratliff | 65 | |
| Bluffport | Sumter | W. H. Sanders | | |
| Bluff Spring | Talladega | M. Street | 90 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 | |

Page Segmentation Result – My Workflow

Page Segmentation Result -- Omnipage

| | | | | |
|----------------|------------|--------------------|-------|---------|
| Blakely | Baldwin | J. H. Stanmyres | | |
| Blake's Ferry | Randolph | J. J. Bradley | | \$17 78 |
| Blocker's | Tuscaloosa | William G. 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. Dendy | | |
| Blue Pond | Cherokee | A. Poore | 20 20 | |
| Blue Spring | Morgan | C. Ratliff | 65 | |
| Bluffport | Sumter | W. H. Sanders | | |
| Bluff Spring | Talladega | M. Street | 90 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 | |
| Breekinridge | Conecuh | W. W. Johnson | | |

| | | | | |
|-----------------|------------|--------------------|--|----------|
| London | Rusk | S. C. Smith | | |
| Lone Star | Titus | W. Collins | | \$123 91 |
| Lone Tree | Collin | John Seaborn | | 3 60 |
| Long Branch | Polina | E. Nix | | |
| Long Point | Washington | M. Rutherford | | |
| Long Prairie | Fayette | R. Gapp | | 12 35 |
| Lookout | Leon | C. J. Dotson | | 2 48 |
| Lynchburgh | Harris | William H. Bryan | | |
| Lyles | Fayette | J. W. Fishburn | | 78 39 |
| McKinney | Collin | E. A. Lowry | | 137 47 |
| McMillan's | Polina | W. McMillan | | |
| Macomb | Grayson | E. McKinley | | 22 08 |
| Madisonville | Madison | J. S. Colvard | | |
| Magnolia | Anderson | William A. H. Hood | | 14 87 |
| Magnolia Spring | Jasper | R. J. Walker | | 18 92 |
| McJomet | Burnet | George Alfer | | 3 23 |
| Malakoff | Henderson | J. A. Mitchano | | |
| Mansfield | arrant | J. Field | | 46 66 |
| Manila | Collin | E. B. Rollins | | 38 70 |

Page Segmentation Result – Adobe Acrobat

Page Segmentation Result – Tesseract

| | | | | |
|----------------|------------|--------------------|--|----------|
| Blake's Ferry | Randolph | J. J. Bradley | | \$17 78 |
| Blake's Ferry | Randolph | J. J. Bradley | | \$817 78 |
| Blocker's | Tuscaloosa | William G. Blocker | | |
| Blocker's | Tuscaloosa | William G. blocker | | |
| Blount Spring | Blount | B. H. Sapp | | |
| Blount Spring | Blount | B. H. Sapp | | |
| Blountsville | Blount | John L. Hopkins | | 66 38 |
| Blountsvfle | Blount | John L. Hopkins | | 66 38 |
| | | H. A. Gillespie | | 33 81 |
| | | H. A. Gillespie | | 33 81 |
| Blowing Spring | Limestone | J. Martin | | |
| Blowing Spring | Limestone | J. Martin | | |
| Blue Lick | Franklin | A. T. Bonds | | |
| Blue Lick | Franklin | A. T. Bonds | | |
| Blue Mountain | Calhoun | J. W. Dendy | | |
| Blue Mountain | Calhoun | J. W. Dendy | | |
| Blue Pond | Cherokee | A. Poore | | 20 20 |
| Blue Pond | Cherokee | A. Poore | | 20 20 |
| Blue Spring | Morgan | C. Ratliff | | 65 |
| Blue Spring | Morgan | C. Ratliff | | 65 |
| Bluffport | Sumter | W. H. Sanders | | |
| Bluffport | Sumter | W. H. Sanders | | |
| Bluff Springs | Talladega | M. Street | | 70 79 |
| Bluff Spring | Talladega | M. Street | | 70 79 |
| Boligee | Greene | S. L. Gwyn | | 8 85 |
| Boligee | Greene | S. L. Gwyn | | 8 85 |
| Bowden | Talladega | J. W. Barton | | |
| Bowden | Talladega | J. W. Barton | | |
| Bradford | Coosa | Thomas Harris | | 25 36 |
| Bradford | Coosa | Thomas Harris | | 25 36 |
| Braggs | Lowndes | R. D. Twombly | | |
| Braggs | Lowndes | R. D. Twombly | | |

Character Error
Rate 1.54%

Table 1. OCR Output – Customized Tool

| 1 | Post Office | County | Postmaster | Compensation |
|----|----------------|------------|-------------------|--------------|
| 2 | Blake's Ferry | Randolph | J J Bradley | \$17 78 |
| 3 | Blocker's | Tuscaloosa | William G Blocker | |
| 4 | Blount Spring | Blount | B. H Sapp | |
| 5 | Blountsville | Blount | John L. Hopkins | 66 38 |
| 6 | | | H. A Gillespie | 33 81 |
| 7 | Blowing Spring | Limestone | J Martin | |
| 8 | Blue Lick | Franklin | A T Bonds | |
| 9 | Blue Mountain. | Calhoun | J W Dendy | |
| 10 | Blue Pond | Cherokee | A Poore | 20 20 |
| 11 | Blue Spring | Morgan. | C Ratliff | 65 |
| 12 | Bluffport | Sumter | W H Sanders | |
| 13 | Bluff Spring | Talladega | M. Street | 79 |
| 14 | Boligee | Greene | 8 L. Gwyn | 8 85 |
| 15 | Bowden | Talladega | J W Barton | |
| 16 | Bradford | Coosa | Thomas Harris | 25 36 |
| 17 | Bragas | Lowndes | R. D Twombly | |

Table 2. OCR output – OmniPage

| 1 | Post Office | County | Postmaster | Compensation |
|----|----------------|-----------------|--------------------|--------------|
| 2 | Blake, Ferry | Randolph | J. J. Bradley | \$17 78 |
| 3 | Blocker, | Tuscaloosa | William G. Blocker | |
| 4 | Blount Spring | Blount | B. H. Sapp | |
| 5 | Blountsville | Blount | John L. Hopkins | 66 38 |
| 6 | | H. A. Gillespie | | 33 81 |
| 7 | Blowing Spring | Limestone | J. Martin. | |
| 8 | Blue Lick | Franklin | A. T. Bonds | |
| 9 | Blue Mountain | Calhoun | J. W. Deady. | |
| 10 | Blue Pond | Cherokee | A. Poore | 20 20 |
| 11 | Blue Spring | Morgan | C. Ratliff | 65 |
| 12 | Bluffport | Sumter | W. H. Sanders | |
| 13 | Bluff Spring | Talladega | M. Street | 20 79 |
| 14 | Boligee | Greene | S. L. Gwyn | 8 85 |
| 15 | Bowden | Talladega | J. W. Barton | |
| 16 | Bradford | Coosa | Thomas Harris | 25 36 |
| 17 | Braggs | Lowndes | R. D. Twombly | 32 26 |

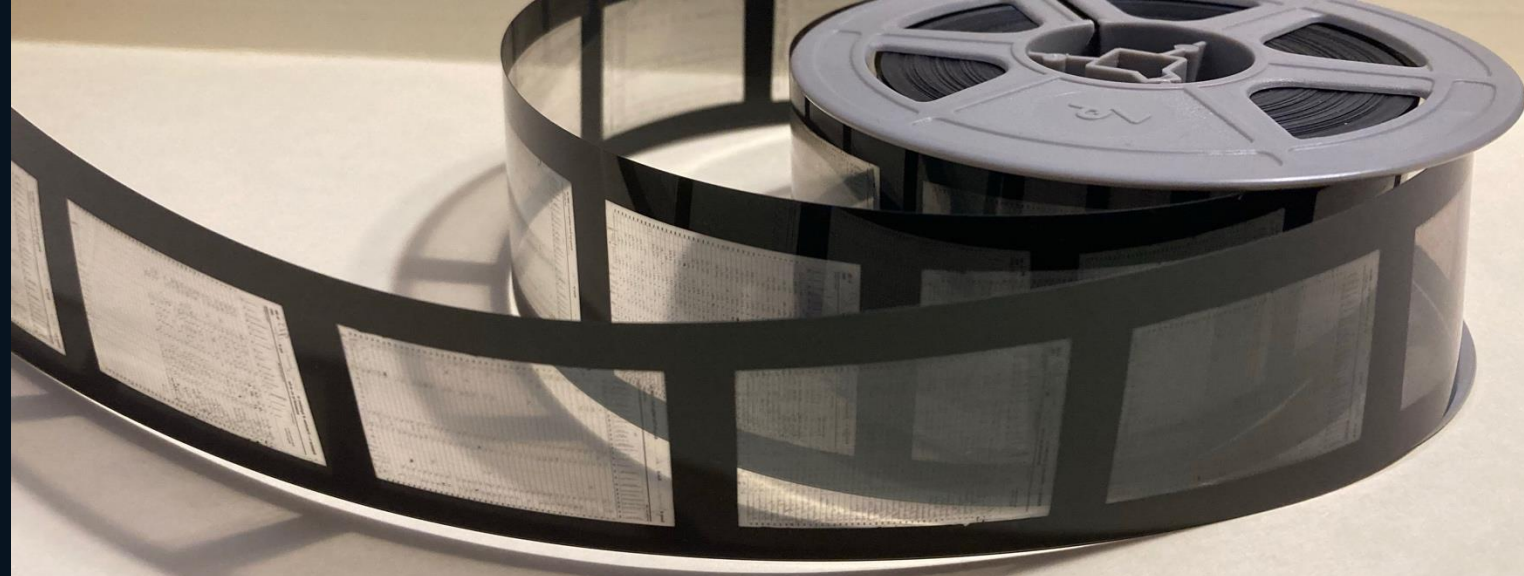
OCR Output ~ Customized Tool v.s. OmniPage

- Page segmentation error by OmniPage
 - Column Splitting Errors
 - Column Insertion Errors
 - Column Shuffling Errors
 - Mixed Errors
- Hard-coded post-processing rules will not work.
 - E.g. “Combine two columns if one is empty.”

Use Case 2:
Reduce
Census
Overcounts

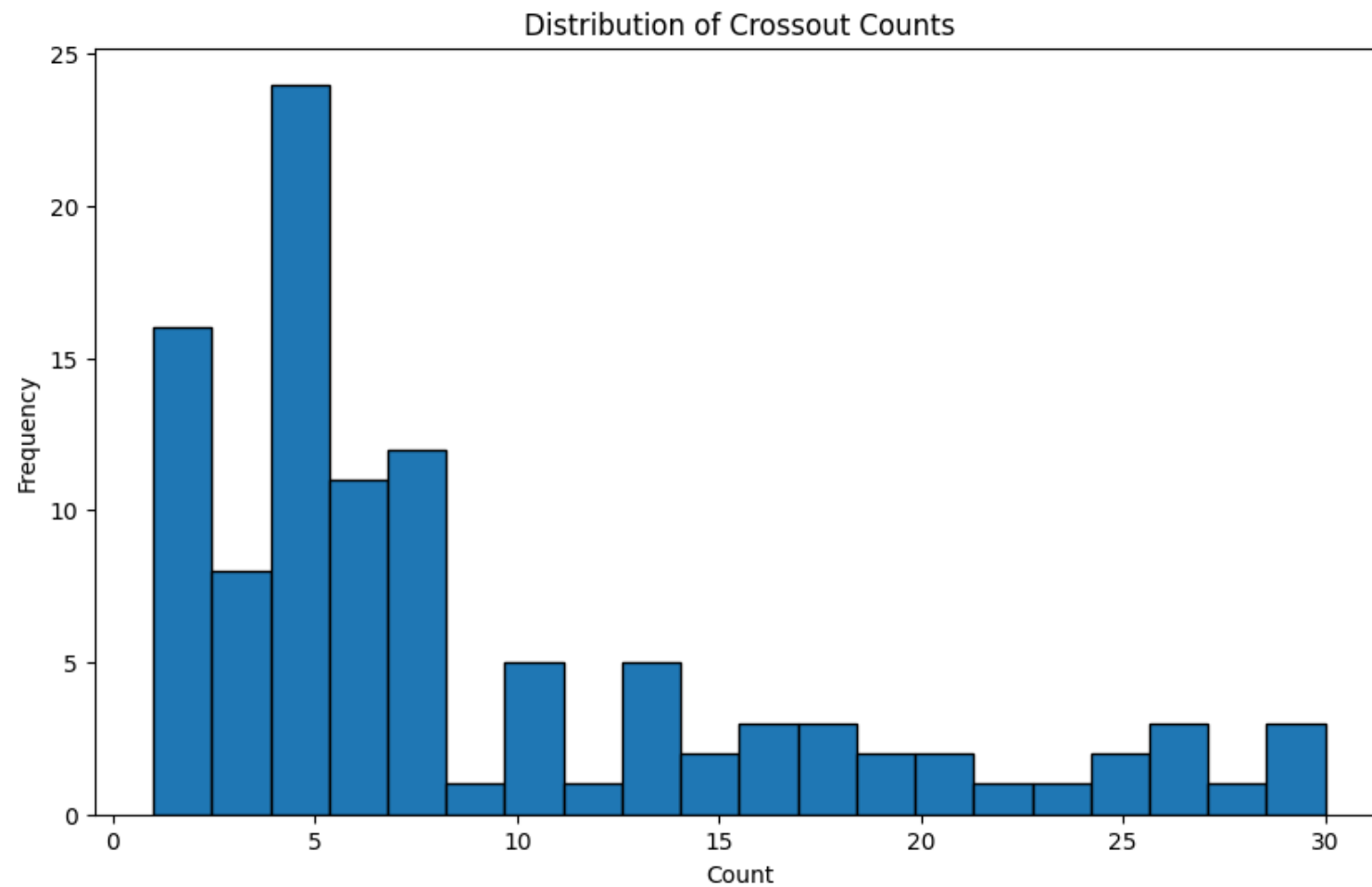


Experiment Design



- 1877 (56310 rows) Census Record pages selected from all states with bias for pages with more strikethroughs.
- Fallback to negatives if an image failed at any stage of the pipeline.
- Calculate validation metrics against manually collected strikethrough information.

Distribution of
Cross-out
Counts in the
Validation
Dataset



Experiment Results

- Precision: 72.64%
- Recall: 91.22%
- **Accuracy: 95.21%**

Confusion Matrix

| | | |
|----------|----------|----------|
| | Positive | Negative |
| Positive | 5705 | 2149 |
| Negative | 549 | 47907 |
| | Positive | Negative |

Actual Class

Predicted Class

- **Precision** measures the accuracy of positive predictions:

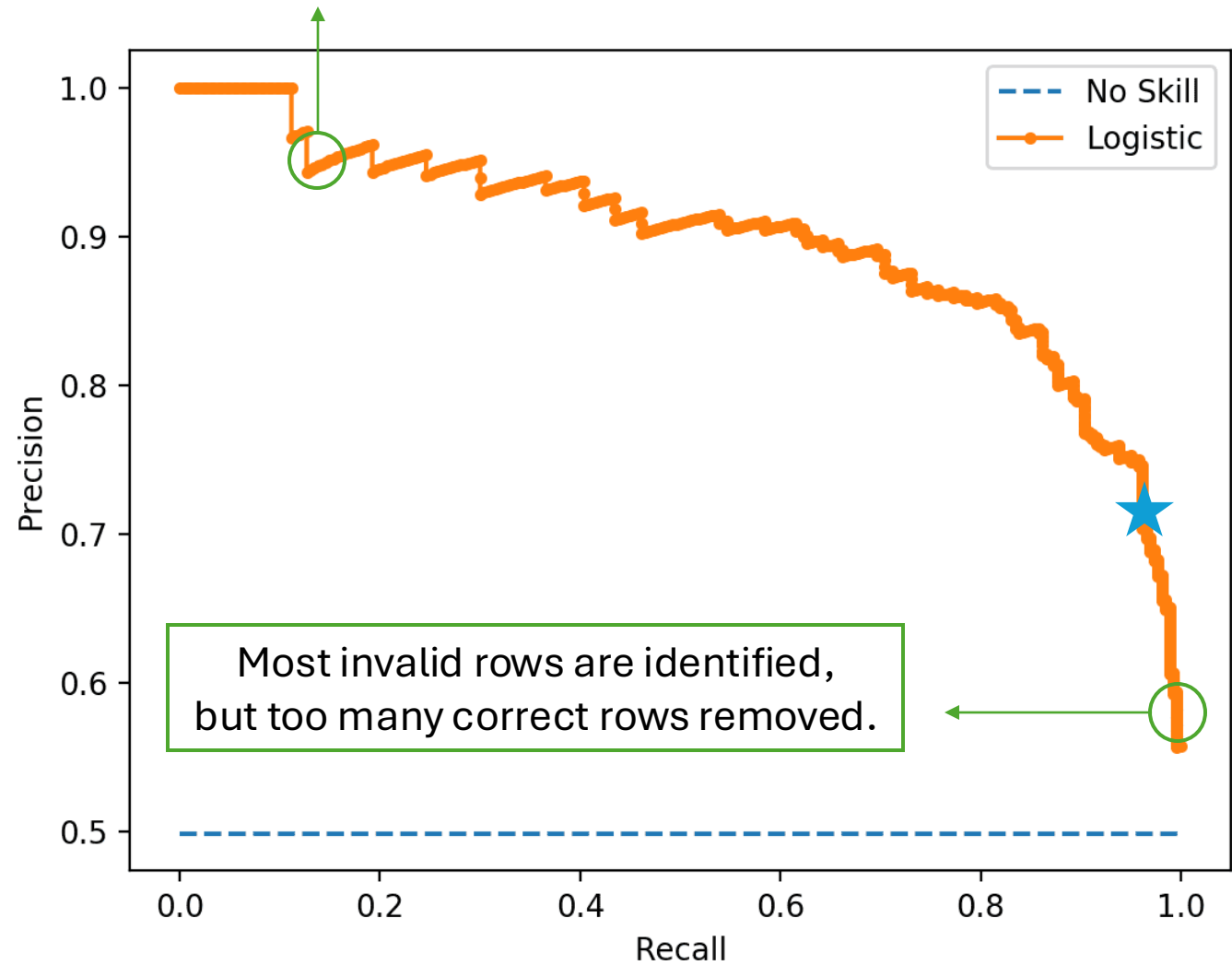
$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}$$

- **Recall** (or Sensitivity) measures the coverage of actual positives:

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}}$$

Precision & Recall Curve

Not enough Rows Identified,
But most rows identified are correct.



Most invalid rows are identified,
but too many correct rows removed.

Thank you!

Email: beczhaozmy0@gmail.com

LinkedIn: <https://www.linkedin.com/in/mengyue-rebecca-z-a15bb8111/>