# 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

MSc. Economic and Social History, University of Oxford, Class of 2019
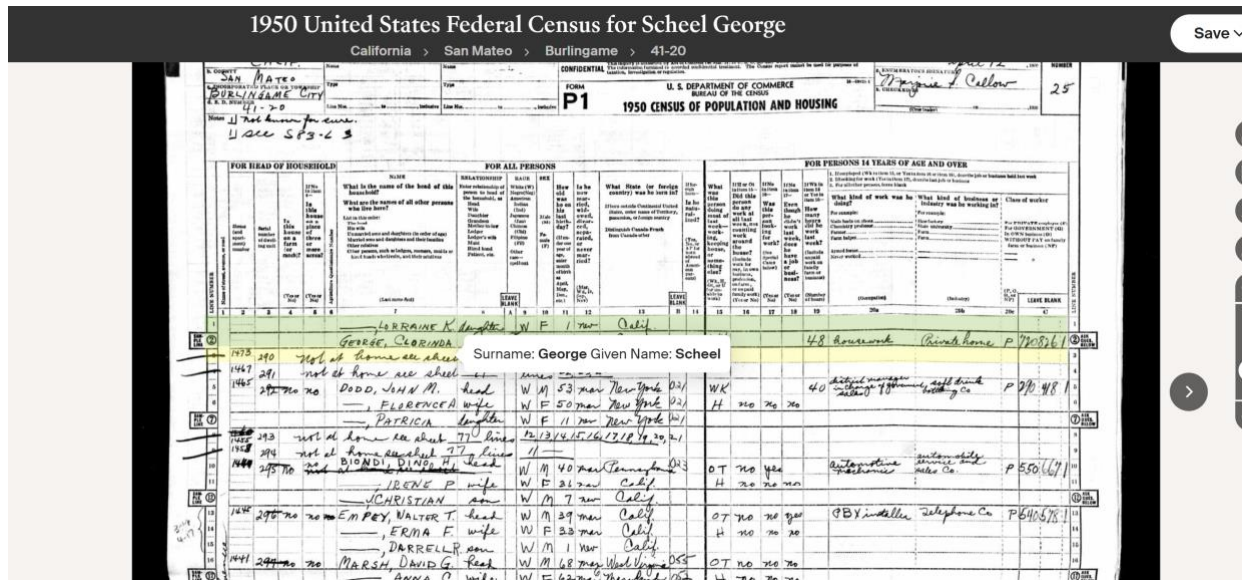
beczhaozmy0@gmail.com

## Abstract

The integrity and reliability of historical records have emerged as critical concerns for researchers in the humanities and social sciences. Traditional methodologies often fall short in addressing the complexities and nuances inherent in historical documents, leading to transcription errors and misinterpretations that can significantly skew data analysis. This study introduces a dual-model deep learning framework that combines semantic segmentation and transfer learning to accurately identify and correct transcription errors across 15 million pages of U.S. Census records. Validated in collaboration with the Minnesota Population Center, the experiment achieved a 95% accuracy rate on a test set of 56,310 entries, effectively eliminating up to 1-2 million erroneous entries and enhancing the trustworthiness of historical data for research applications.

## Introduction

This study applies cutting-edge deep learning and AI models to improve data accuracy across 15 million pages of U.S. Census records by identifying and correcting transcription errors and erroneous entries. Through a human-validated test set of 56,310 entries, the project achieved 95% accuracy in correcting transcription errors, eliminating an estimated 1-2 million errors. This sets a new benchmark for integrating technology into social sciences, with broad implications for improving research reliability in any endeavor requiring large-scale data collection from historical sources.

## Problem Statement

The 1950 census shows an estimated 1-2 million overcounted entries across 15 million pages when cross-referenced with county-level data. Traditional OCR methods struggle with these inaccuracies, such as confusing "Noah" with "No One." My project, conducted in partnership with the Minnesota Population Center, utilizes deep learning to address these inaccuracies, achieving a 95% accuracy rate in error detection and correction, ultimately enhancing the reliability of historical data for research. In the example bellow, row 3 says "not at home see sheet [...]" but Ancestry.com detects "George Scheel" based on visual similarity.

## Limitations of Traditional Computer Vision for Complex Historical Documents

Edge detection and traditional computer vision methods are poorly suited for processing the kind of complex information seen in historical census documents like those in the image. These records contain dense, overlapping elements, such as strikethroughs, circles, and crossed-out numbers, which would confuse an edge detection algorithm. Such algorithms rely on detecting sharp contrasts and well-defined lines, but the handwritten marks and varied structures in these documents do not fit such a pattern. Overlapping strokes and faded lines are common, further complicating the task. The inconsistent handwriting, varied thickness of lines, and degradation artifacts introduce noise that would likely be misinterpreted as valid edges or boundaries, producing inaccurate results. Without a semantic understanding of the content, edge detection would fail to distinguish between meaningful data points and irrelevant markings, leading to errors and missed information.

Moreover, the structure of historical census records lacks the standardization that edge detection requires. Each entry may contain unique notations, symbols, and inconsistent row structures, making it challenging for basic algorithms to generalize. Edge detection alone would struggle to parse these subtle distinctions, as it lacks the ability to classify elements within the context of a historical document's layout. In contrast, a deep learning framework equipped with semantic segmentation and image classification can be trained to recognize and interpret these complex patterns and contextual nuances, enabling it to differentiate between data points, noise, and layout elements. This approach

allows for a higher accuracy rate, as it can adapt to the specific irregularities in each record, handling the variability in handwriting, formatting, and document quality that edge detection cannot manage effectively.
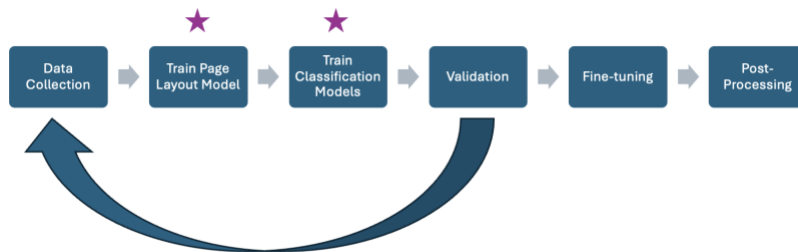
A subset of examples accurately identified by the Deep Learning Framework:



# Dual Deep Learning AI Framework for Precision Data Correction

The Dual Deep Learning AI Framework, honored with the Digital Innovation Award at the 2022 World Economic History Congress, integrates two specialized models: a page segmentation model and a classification vision model/OCR engine. This dual-model framework extracts and corrects data from large, unstructured historical datasets like census records with high precision, offering unprecedented reliability for research dependent on vast primary sources.

Dual Deep-learning AI framework Iteration Process



This diagram outlines the iterative process of developing a machine learning framework for historical document analysis, showing the key stages involved from data collection to post-processing. Here's a breakdown of each step:

1. Data Collection: Gather a large set of census images with accurate labels, prioritizing pages with strikethroughs for comprehensive training.

2. Train Page Layout Model: Train a model to segment uniform census structures, focusing on rows, columns, and marked areas for downstream classification.

3. Train Classification Models: Develop models to identify specific content and patterns, particularly strikethroughs, within segmented areas.

4. Validation: Validate using a test set of 1,877 pages with human-verified data, with a bias in favor of pages with a lot of strike-throughs.

5. Fine-Tuning: Refine models by focusing on challenging or special cases.

6. Post-Processing: Build a pipeline to format outputs into human-readable data, ensuring clean, usable results.

7. Iteration (Feedback Loop): Continuously improve by feeding insights from post-processing back into data collection and training.

## Intelligent Page Segmentation Model

Page Segmentation is a critical yet often overlooked step when working with historical documents. While commercial or open-source OCR software performs well on modern, clearly structured tables, it struggles with the unique complexities of historical formats. The example below demonstrates the advantage of a customized solution designed for historical data, like the *Official Register of the United States*. Here, OmniPage fails to accurately detect rows due to the dotted line structure, which renders the resulting data unreliable for downstream use.



Figure 2. ROI Detection Results -- Customized OCR Tool



Figure 3. ROI Detection Results -- OmniPage

## Ensuring Trustworthy Results with Precision–Recall Balance

Balancing error detection and valid entry retention is essential in data validation. In this census project, preserving valid data was prioritized, even if a few erroneous entries were retained. This approach maintains high data integrity, as demonstrated by the model's confusion matrix. Despite the dataset's imbalance, the model reliably identifies both valid and erroneous entries, supporting the project's goal of maintaining a high standard of data trustworthiness.

The model strikes a careful balance between retaining valid data (negative class) and detecting errors (positive class). With a high specificity of 95.71%, it effectively keeps important records intact, supporting the project's data integrity goals. Additionally, the model's recall rate of 91.22% ensures it captures most errors accurately, while a precision of 72.63% indicates some tolerance for false positives, which can be easily removed by downstream processes. Overall, the model maintains a reliable dataset by preserving most valid entries while effectively identifying errors.

## Confusion Matrix



Confusion Matrix

|  | Predicted Positive | Predicted Negative |
|---|---|---|
| **Actual Positive** | 5705 | 2149 |
| **Actual Negative** | 549 | 47907 |

**Generative AI Model vs. Deep Learning Framework: Performance Comparison**

**Experiment Design:**

This experiment tested the Gemini-1.5-pro model's ability to detect crossed-out line numbers in 1950 U.S. Census forms. Gemini-1.5-pro is one of the most performant models for multi-modal tasks on the market, it provides good baseline understanding of Large Language Models' performance on such tasks. Using a few-shot example approach, the model was tasked with identifying visibly crossed-out rows across 10 pages, where rows are marked by a line, stroke, or scribble through the line number.

***Prompt Details:***

"Context: The image shows a section of a 1950 U.S. Census form, with line numbers listed in the far-left column. Rows may be crossed-out to indicate corrections, removals, or other changes in the census data.

Task: Identify and return a list of line numbers that are visibly crossed-out in the provided image of the 1950 U.S. Census form. A row is considered crossed-out if there is a mark, such as a line, stroke, or scribble made with a pen, pencil, or marker, striking through or across the entire line number in the leftmost column."

***Example Image:***

See Appendix A.

Example LLM Output:

10, 11

***Analysis of Output:***

The actual strikethrough spans lines 8 to 30, meaning all rows within this range should have been marked as crossed-out. By only identifying lines 10 and 11, the model missed a substantial portion of crossed-out lines, indicating LLM's limitations in detecting extended patterns and composite layouts in historical records.

***Summary:***

The experiment, which included a few-shot example approach and testing across 10 pages, suggests that while Gemini-1.5-pro might perform better with additional few-shot examples or fine-tuning, its current output does not look promising for complex pattern recognition in historical document layouts. The model's limitations in handling composite page-layout and pattern analysis tasks are evident. In contrast, my deep learning framework, which integrates semantic segmentation and contextual analysis, is better suited to accurately detecting intricate patterns and preserving high data integrity within archival records.

## Conclusion

This research highlights the transformative potential deep learning in advancing social sciences and humanities research. By leveraging advanced AI for historical document analysis, this study significantly improves the reliability of historical records, contributing to broader discussions on identity, kinship, and community in the digital age. The project establishes new benchmarks for integrating AI in historical research, deepening our understanding of historical records' relevance to modern analysis. While multi-modal LLMs show promise, they are not yet a replacement for conventional deep-learning vision techniques, as they struggle with composite page layouts and complex pattern analysis.

# Bibliography

Ancestry.com. *1950 United States Federal Census*. Lehi, UT: Ancestry.com Operations, Inc., 2022.

# Appendix A



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

State: MASS
County: ESSEX
Incorporated place or township: AMESBURY
E.D. No.: 5-6
Date sheet started: MAY 26
Sheet Number: 80

| Line | House number | Name | Relationship | Race | Sex | Age | Marital status | State or country of birth | Occupation | Industry | Class of worker |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 265 | VIGEANT, MARIE R. | HEAD | W F | 79 | Wd | Canada French | | | |
| 2 | 266 | OBERT, NELLIE L. | HEAD | W F | 55 | Wd | New Hampshire | Bench worker | Shoe factory | P |
| 3 | 265 | OBERT, NORMAN E. | HEAD | W M | 26 | Mar | New Hampshire | Auto Mechanic | Filling Station | P |
| 4 | | JEAN M. | WIFE | W M | 26 | Mar | New Jersey | Bench worker | Electrical Plant | P |
| 5 | 266 | LeBlanc, Franklin | Head | W M | 25 | Mar | Mass. | Machinist | Mfg Shoe Fdy | |
| 6 | | Juliette | Wife | W F | 24 | Mar | Mass. | Bench worker | Mfg Shoe Fdy | |
| 7 | | Blouin, Paul W. | Brother-in-law | W M | 22 | Nev | Mass. | Steeple Jack | Chimney Co. | P |

Persons transcribed from I.R.
Persons not assigned to dwelling unit

Sawyer, John
Powell, Everett C.
Osborne, Ruth E.
Reid, Clara

THE QUESTIONS BELOW ARE FOR PERSONS LISTED ON SAMPLE LINES
FOR ALL AGES / FOR PERSONS 14 YEARS OF AGE AND OVER