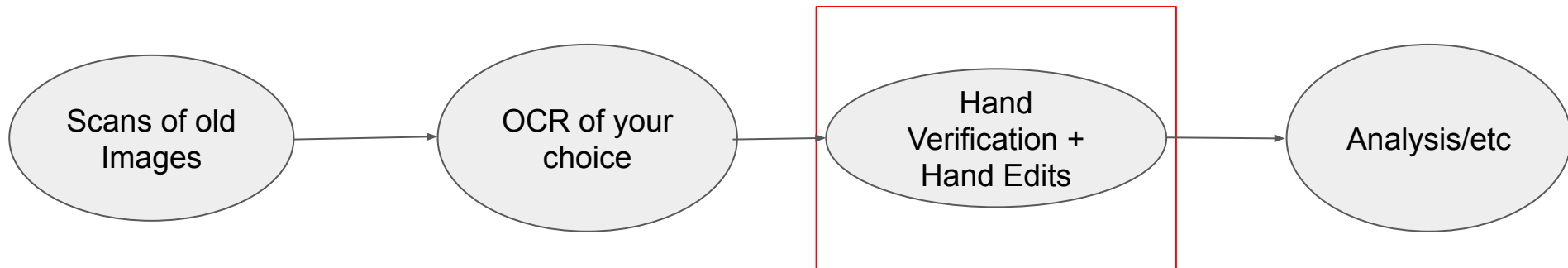# Record Linkage GUI

Todd Nobles, Natalie Turner, Terresa Tran, Honglam Van, Julia Zhu

# Background

- Large scale quantitative analysis of archival data involves a number of steps that are tedious, error prone, and create difficulties for reproducibility
- Our tool assists with the first two of these problems by improving the workflow for users who are turning images of historical documents into analysis ready datasets

```
Scans of old
   Images        →      OCR of your
                          choice        →     Hand
                                           Verification +
                                            Hand Edits     →    Analysis/etc
```

# Data Used

We use two data sources:

1. Images of archival documents
2. A CSV of the computer extracted information from these images

# Use Cases

- Researchers working on **digitizing** and **transcribing** historical forms
    - Compare side by side text fields extracted from an image and the original image for accuracy
    - Manually edit the text for a field that was extracted from an image and save a new dataset with the corrected data
- Researchers with **little to no experience in coding**
    - Simple interface that allows user to quickly check information extracted from images without needing to work with file paths or code

# Demo

# Design

Core components:

- User file upload
- Text display
- Image display
- Editable text fields

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│  Upload image(s)│ ───▶ │     Compare     │ ───▶ │ Make edits to   │ ───▶ │    Save out     │
│  and extracted  │      │ displayed image │      │     text        │      │  updated csv    │
│   text csv file │      │   text fields   │      │     fields      │      │      file       │
└─────────────────┘      └─────────────────┘      └─────────────────┘      └─────────────────┘
```

- record_linkage_gui
  - |- README.md
  - |- LICENSE
  - |- environment.yml
  - |- pyproject.toml
  - |- src
    - |- record_linkage
      - |- app.py
      - |- init.py
  - |-tests
    - |- test_app.py
    - |- test_data
      - |- sample_data.csv
  - |- docs
    - |- user_story.md
    - |- use_cases.md
    - |- components.md
    - |- Record Linkage Demo.mp4
  - |- .github
    - |- workflows
      - |- testsuite.yml
  - |- .gitignore

# Repository Structure

# Lessons Learned and Future Work

- Sometimes documentation of specific components of a UI framework is not sufficient for newcomers to the tool. Full working examples are more helpful
- Our timing was off: we moved to testing in class before we had a UI, so we had to rush to build something testable.
- Our testing didn't always fit the standard testing discussed in class
  - We tested underlying functions with pytest, but to test Streamlit's UI behavior we did manual verification
- **Future work**
  - Improving scalability, currently a user would be limited to as many images as their RAM can store in working memory
  - Building the second tool in the Record Linkage Toolkit to help users adjudicate between potential matches across datasets by hand
  - Further testing once Streamlit has more capabilities to mock user inputs and interactions with the application