**THE UNIVERSITY OF CHICAGO**

# MARI:
## Semi-Automated, Human-in-the-Loop Redaction of Text Corpora
Emma I. C. Peterson, Valerie Zhao, Dan Byrne, Blase Ur
University of Chicago

SUPER GROUP

## Motivation

- Sharing qualitative, naturalistic data for the benefit of multiple research groups is difficult
- Protecting participants: redacting identifiable information
  - Most existing *automatic* redaction tools are insufficient
    - Balancing privacy (redacting information) vs research contribution (leaving data untouched)
  - Manual redaction is impractical due to budget and time constraints

## Contributions

- **MARI (Mostly Automated Redaction of Identifiers),** a prototype **human-in-the-loop** tool informed by interviews with social science researchers
- Interview findings on redaction and research methods
- Unique combination of PII feature engineering, linguistic analysis, and information-theoretic scoring
  - *Goal:* maximize data utility, cross-discipline generalization, comprehensive redaction coverage

## Taxonomy and Examples

| Taxonomy Category | Taxonomy Sub-Categories | Examples |
|---|---|---|
| Identifiers | Personal Names, Nicknames, Personal Identity, Numbers | *Legally given name, Diminutives, SSN, EIN* |
| Demographics | Age, Sex, Gender, Pronouns, Sexuality, Race & Ethnicity, Education, Profession, Health Status | *Date of Birth, LQBTQ+, Niche Job, Rare Disease, HIV Positive, Mixed Race* |
| Locations | Country, State, City, Postal Code, Address, Landmark, Business | *United States, Illinois, Chicago, 5307 S Woodlawn Ave, The Bean, Jimmy's Tap* |
| Dates & Events | Publicly Recognized Dates & Events, Personal Dates & Events | *Thanksgiving, Christmas, Cancer-Remission Anniversary* |
| Linguistic Patterns | Regional Dialects, Code-Switching, Unique Vocabularies, Idiosyncratic Expressions | *African American Vernacular English, Scots-English, Spanglish, Parmesan Cheese == "Pasta Sugar"* |
| Personal Interests & Activities | Traditions, Group Membership, Cultural References, Popular Culture Participation, Hobbies | *University / school traditions, belonging to native tribe or military, member of a small fandom* |

*Table 1*. Our privacy taxonomy consists of hierarchical categories and highlights certain redaction decision thresholds. Classifiers will be constructed per category and tuned to different data presentation types.
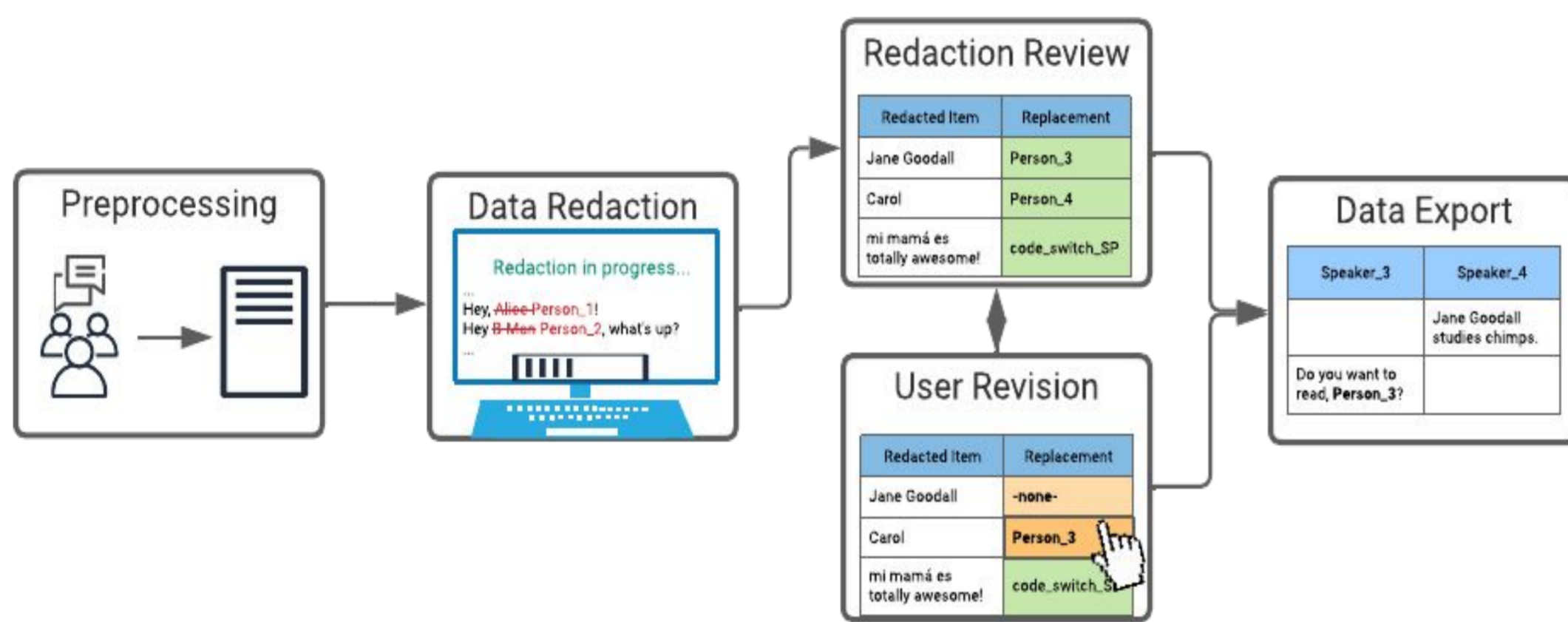
## System Implementation



*Figure 1*. Application flow diagram. The user uploads data, which is then redacted by our system; the user can then revise the automated decisions.
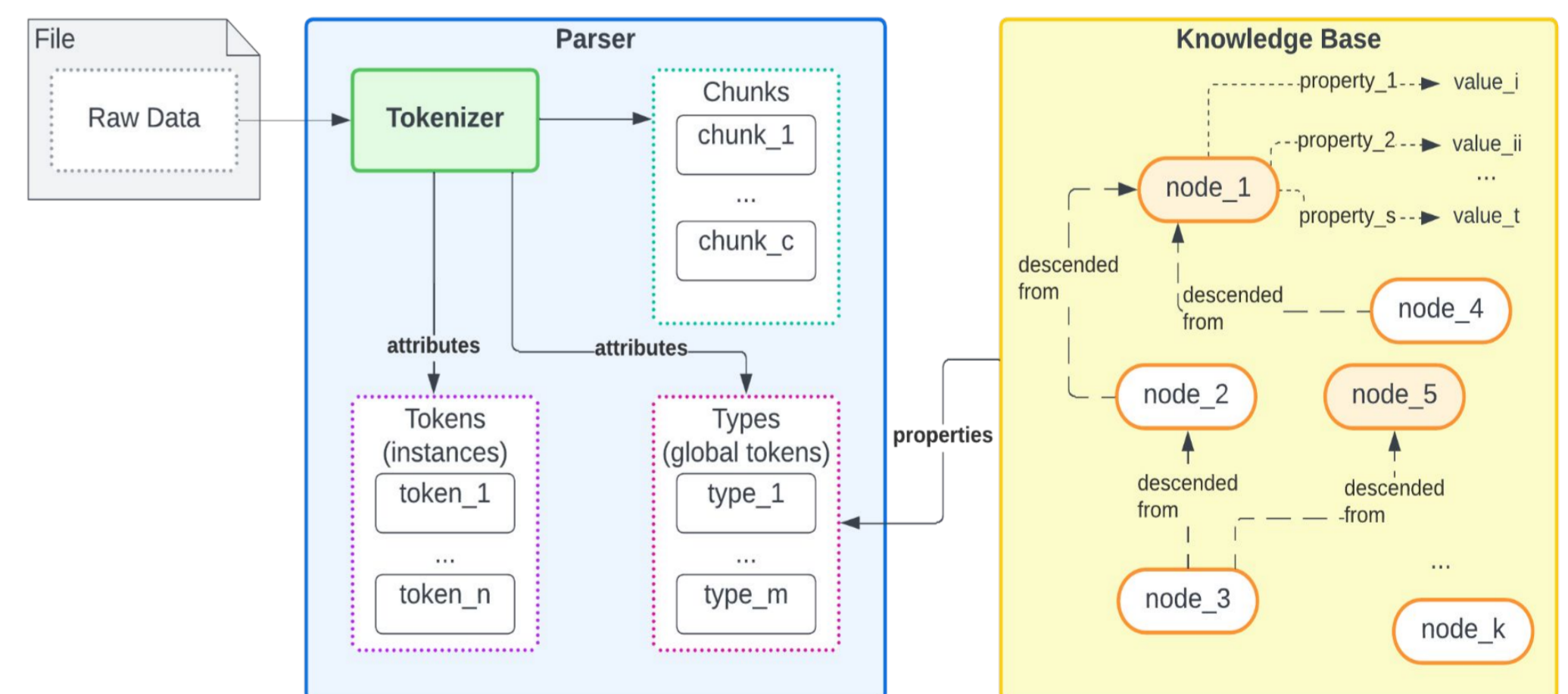


*Figure 2*. Overview of application pre-redaction processing and connection of types to knowledge base.
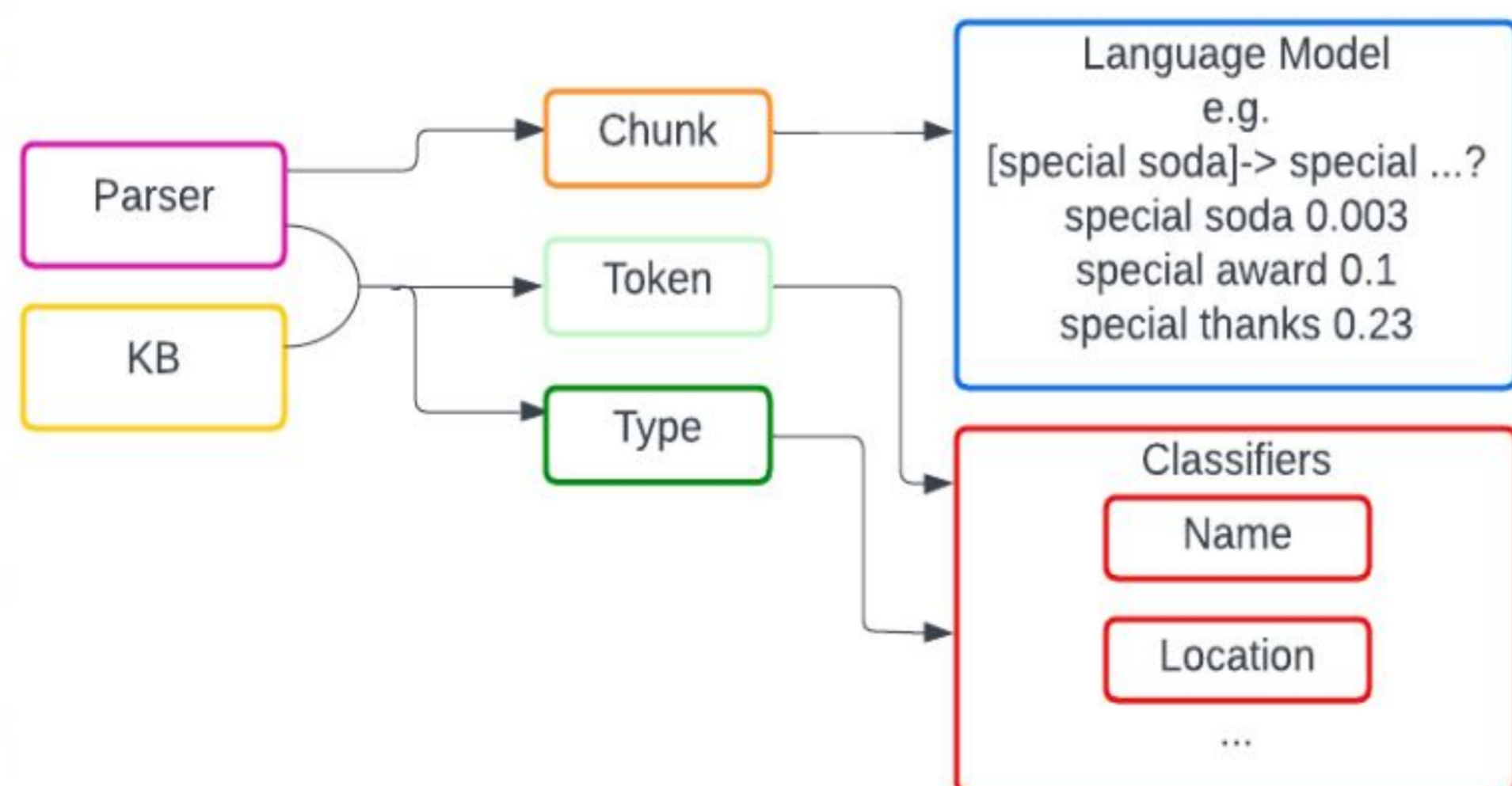


*Figure 3*. Redaction flow diagram. Data is parsed into different structures. Tokens and types are given to classifiers, whereas chunks are analyzed by language modeling.
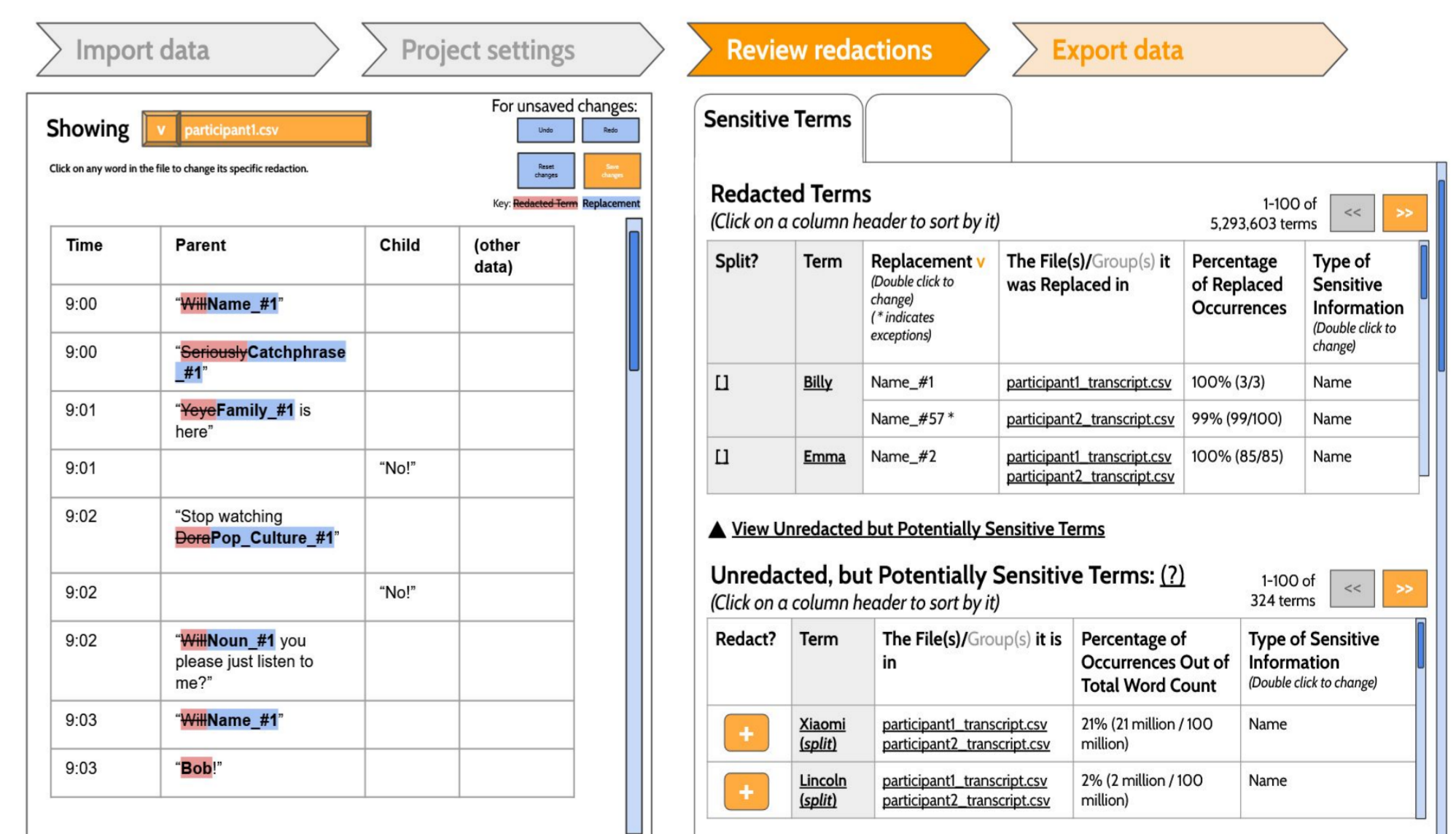


*Figure 4*. UI mockup displaying redacted information to the user. Redacted terms are shown to the user in context both file-wise and term-wise.