

Detecting Informal Data Use in Literature

Sara Lafia¹, Elizabeth Moss¹, Andrea Thomer², and Libby Hemphill^{1,2} (University of Michigan: ICPSR¹; UMSI²)

Challenge

Formal data citations using unique identifiers are readily discoverable; however, informal references indicating research data reuse are challenging to detect. *How can computational approaches to detect data use complement human efforts?*

Procedure

- Search for formal (unique identifiers) and informal mentions (study names, aliases) of research data
- Extract terms (“survey”, “sample”...) that often accompany data citations and sections of articles (Methods...) where found
- Predict custom entity type (*Data*) at the sentence level
- Evaluate candidate documents for inclusion in the ICPSR Bibliography of Data-Related Publications

Detecting citations: human vs. computational approach

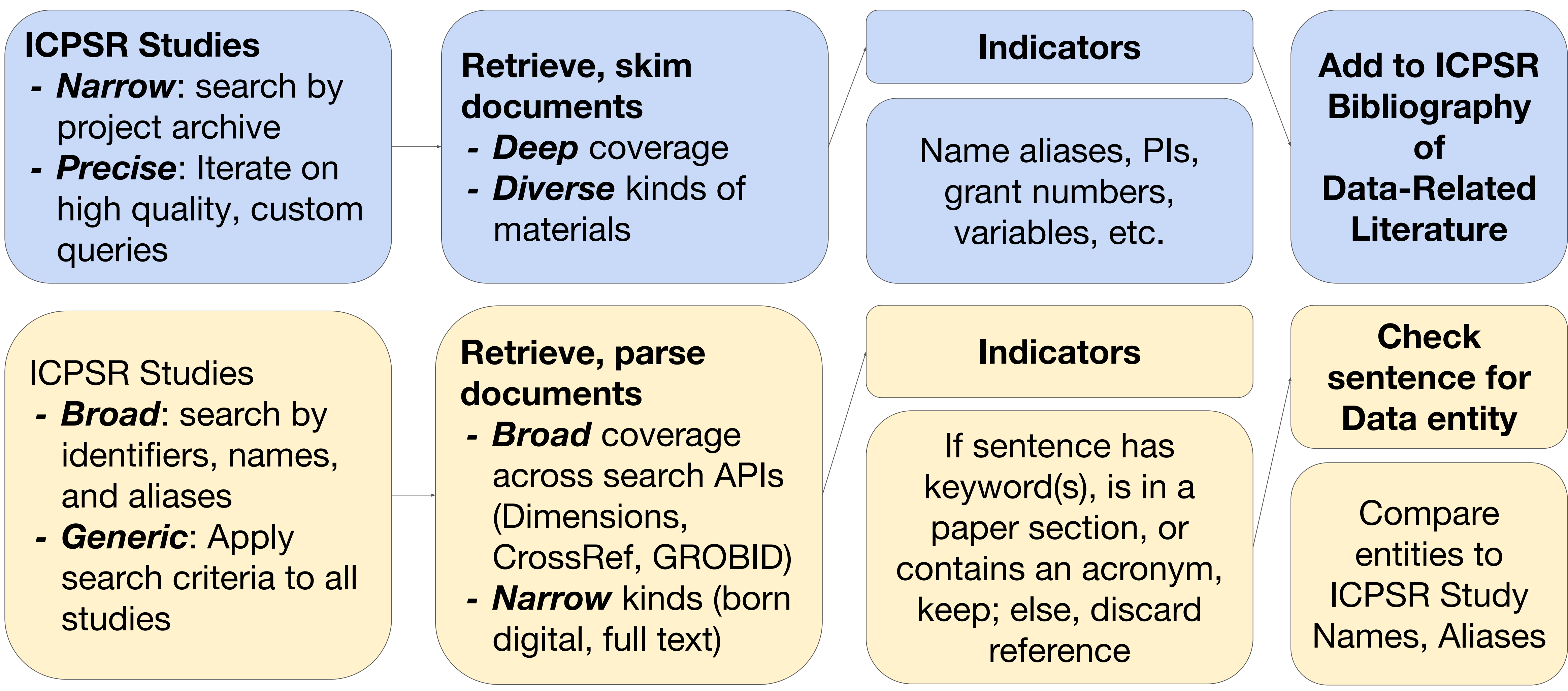


Figure 1. Comparison of human heuristics (left) and a computational approach (right) using a custom Named Entity Recognition model trained to predict passages of text indicating research data applied to the article, **The political legacy of American slavery** (Acharya et al., 2016). doi: [10.1086/686631](https://doi.org/10.1086/686631)

Left (Human Heuristics):

Political Legacy of American Slavery Acharya, Acharya, Matthew Blackwell, and Maya Sen

should be weaker (i.e., should have decayed more) in areas where the incentives for anti-black attitudes faded earlier.

How contemporary demographics could explain regional variation in white political attitudes

In contrast to the arguments above, much of the political science literature points to contemporary (not historical) forces as providing the explanation for why Black Belt whites are more conservative on race. By and large, the literature has interpreted Key's (1949) work as suggesting that whites contemporaneously become more conservative when they are exposed to the high concentrations of African Americans who live in their communities. The high concentration of African Americans in today's Black Belt could contemporaneously threaten white dominance, resulting in whites actively choosing more conservative political beliefs today. The literature supporting this idea, known as "racial threat," is voluminous. For example, Glaser (1994) finds evidence linking negative white attitudes toward civil rights or African American politicians with high concentrations of blacks. Glesne and Brinkner (1993) find a relationship between black concentrations and white support for racially conservative candidates such as David Duke (these findings are, however, challenged by Voss [1996]). This literature, however, has not considered that slavery could be an independent predictor of contemporary attitudes (apart from its effect on contemporary demographics), making it an omitted variable in studies of racial threat in the South.

Other aspects of the contemporary local context may also affect white attitudes—for example, income gaps between blacks and whites, urban-rural differences, and other contextual and individual-level factors (e.g., Hopkins 2010; Oliver and Mondak 2000). A final category of explanations concerns white mobility through the twentieth century. For example, it could be that more racially conservative whites have migrated into former desegregating areas, while racial liberals have left, thereby creating a regional pattern in the South.

Right (Computational Approach):

Since county boundaries have shifted since 1860, we use an area-weighting method to map data from the 1860 Census onto county boundaries in 2000 **DATASET**, enabling us to estimate the proportion enslaved in 1860 within modern-day counties.

We analyze three county-level outcome measures, which come from the Cooperative Congressional Election Study (CCES **DATASET**), a large survey of American adults (Ansolabehere, 2010).

We pool CCES **DATASET** data from the 2006, 2008, 2009, 2010, and 2011 surveys to create a combined data set of over 157,000 respondents.

In addition, we also investigate individual-level black-white thermometer scores from waves of the American National Election Survey (ANES **DATASET**) from 1984 until 1998, a time period where the ANES both used a consistent sampling frame and included county-level identifiers for respondents.

After restricting the sample to Southern whites, we have an ANES **DATASET** sample of 3,123 individuals across 64 counties in the South.

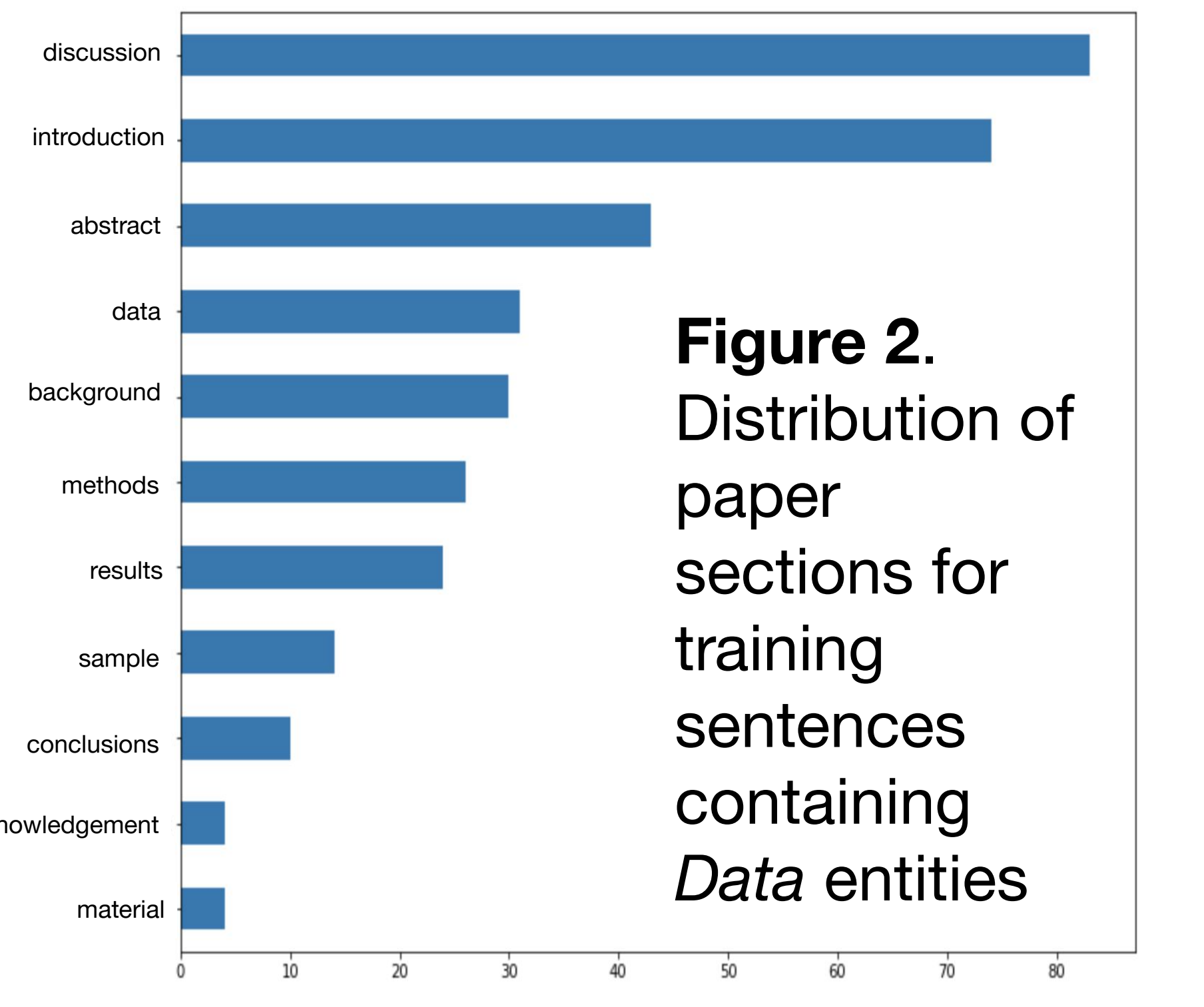
This makes the ANES **DATASET** more restricted in its geographic coverage, but it contains valuable direct questions on the subjective evaluation of racial groups.

We construct our partisanship measure from a standard seven-point party identification question on the CCES **DATASET**.

All CCES **DATASET** surveys ask respondents whether they support or oppose affirmative action policies, which are described as "programs [that] give preference to racial minorities and to women in employment and college admissions in order to correct for discrimination" (2008 CCES).

Training a computational model

1. Label *Data* entities in 2,056 sentences from 400 papers
2. Analyze distribution of indicator terms by section
3. Train custom spaCy NER pipeline with annotated sentences



Acknowledgements

This material is based upon work supported by the National Science Foundation under grant 1930645.

