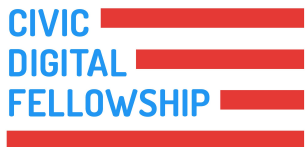# CHARACTERIZING THE NIH RESEARCH WORKFORCE

## NIH Office of the Director, Office of Extramural Research

Katrina Pearson, Director of Division of Statistical Analysis and Reporting

CIVIC DIGITAL FELLOWSHIP

National Institutes of Health
Office of Extramural Research

SAM KIM
Duke University
Biology & Computer Science

# MOTIVATION

- Annual NIH appropriation exceeds $37 billion*

- More than 80% is awarded to researchers as **extramural grants/contracts**[†]

- Research output can be inferred from publications

- What is the **workforce output**?

  - That is, **how many jobs** does NIH funding directly create?

CIVIC
DIGITAL
FELLOWSHIP

NIH

# MOTIVATION (cont.)

- **Automate and generalize** the workforce analysis



**Analysis of FY 2009**
(Pool et al., 2016)*

**Analysis of FY 2017**
(Kim et al., 2019)

* Pool LR, Wagner RM, Scott LL, et al. (2016). Size and characteristics of the biomedical research workforce associated with U.S. National Institutes of Health extramural grants. *FASEB J.* 30(3):1023–1036.

# DATA SELECTION

- **All-Personnel Reports** (APRs)

  - Required for RPPRs since 2010

  - Self-reported by awardees

  - Includes all people who devoted 1+ month of effort

- **All awarded and funded NIH extramural grants**

  - No intramural grants or inter-/intra-agency agreements

  - No subprojects or contracts

  - No noncompeting supplements (Type 3's)



IMPAC II
Information for Management Planning Analysis and Coordination

OLTP
Online Transaction Processing Database

IRDB
IMPAC II Reporting Database

QVR
Query, View, and Report System

# DATA PROCESSING

- **Very poor data quality**

- Misspelled names/words

- Position titles in free text

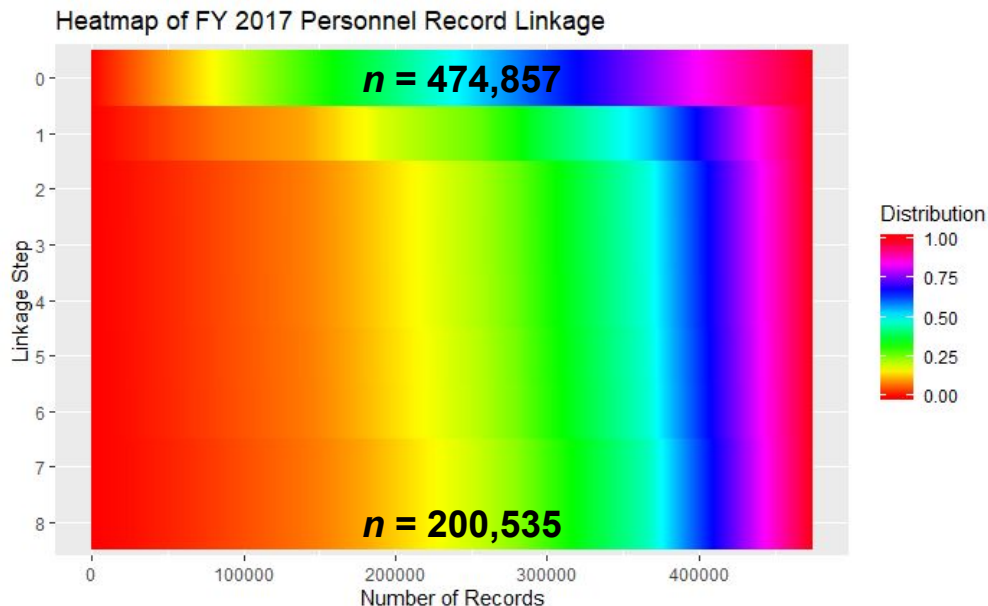    - Custom spellchecker using edit distances

- Missing identifiable information

    - Record linkage methods

- **Missing RPPRs** for about 28% of grants

    - Impute personnel data



Distribution of word frequencies
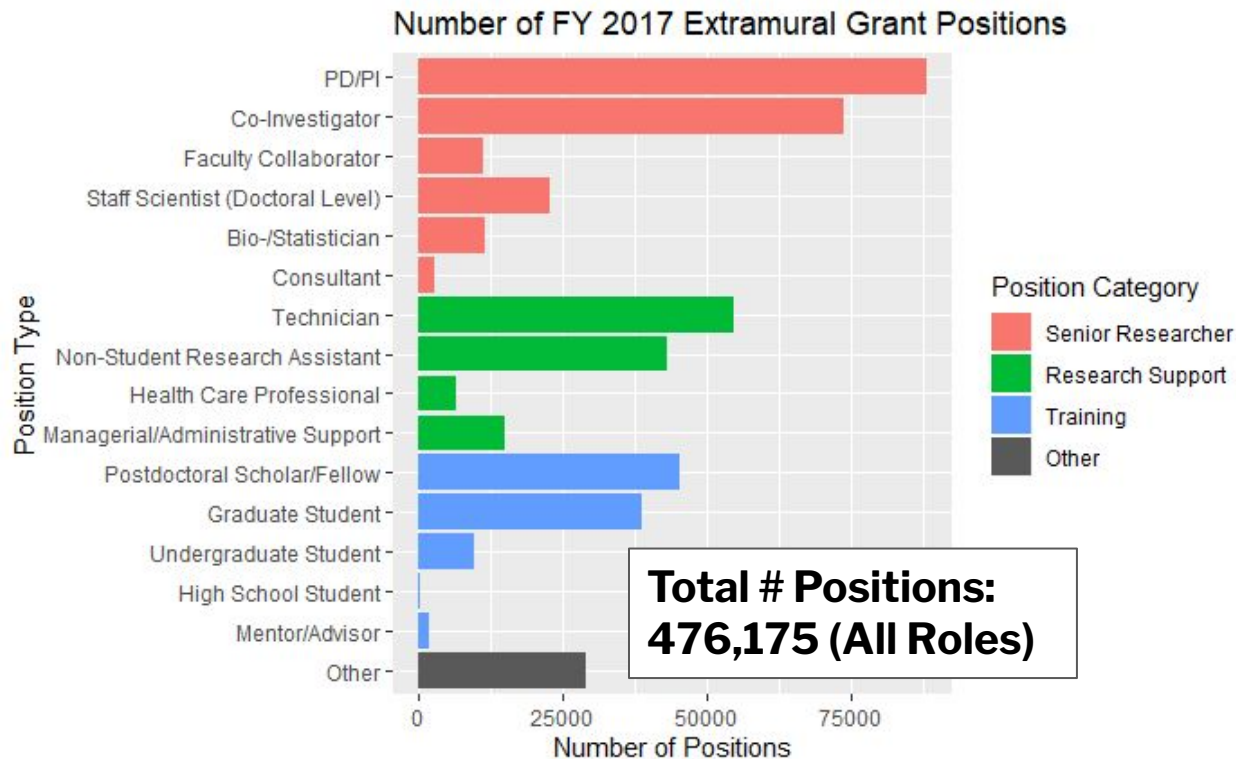
**96% coverage
(Top 177 words)**

# RECORD LINKAGE

- How to resolve duplicated, misspelled or similar names?

    - First and last names (99%)

    - Middle names (33%)

    - Commons profile ID (58%)

    - Institutional ID (100%)

    - SSN, birthdate (13%)

- Machine learning was impractical

    - Clustering: 450k × 450k

Heatmap of FY 2017 Personnel Record Linkage

$n$ = 474,857

$n$ = 200,535

Linkage Step

Number of Records

Distribution
1.00
0.75
0.50
0.25
0.00

# # POSITIONS



Number of FY 2017 Extramural Grant Positions

Total # Positions: 476,175 (All Roles)

# # POSITIONS



Number of FY 2017 Extramural Grants

Average Number of Positions Per Project

# # UNIQUE PEOPLE



Number of FY 2017 Extramural Grant Personnel

**Total # People:**
**200,535 (Pre-Imputation)**
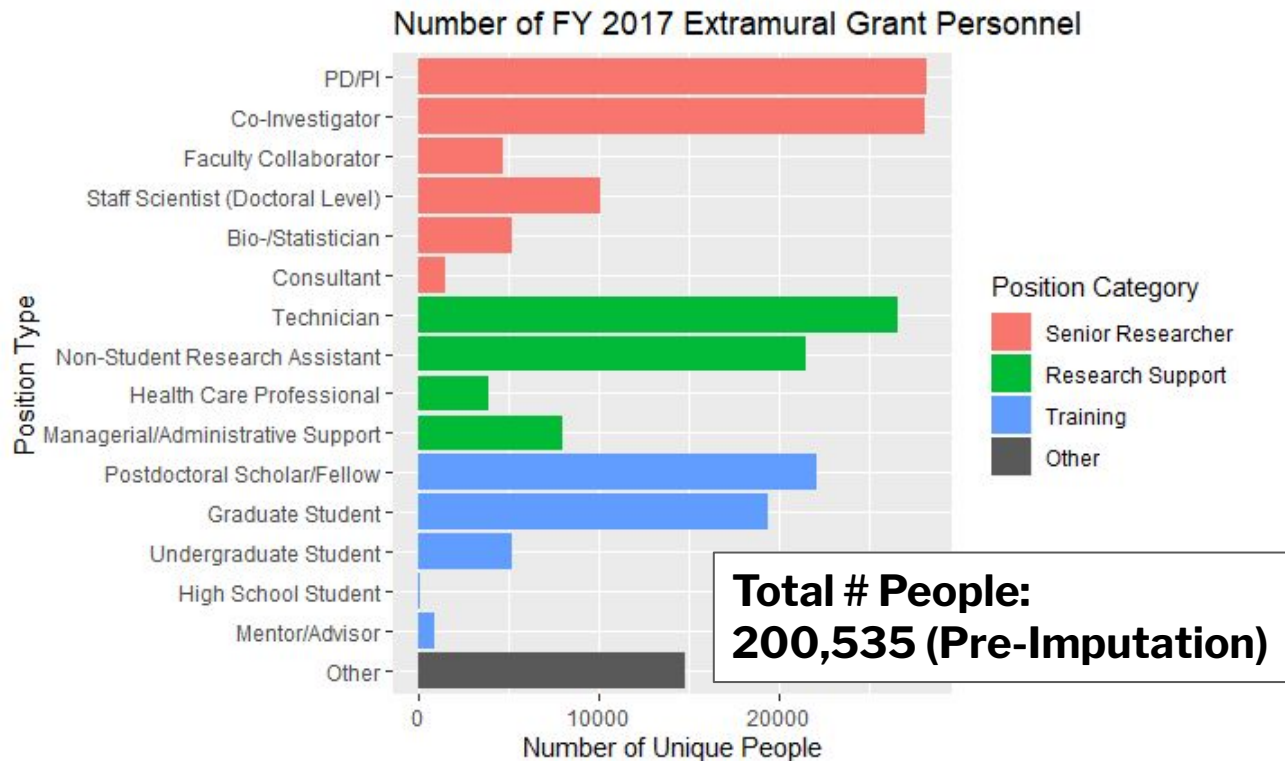
# NEXT STEPS

- Impute missing personnel data, especially for Type 1's

- Breakdown by geographic location

- Analyze effort and career stage information

- Analyze for FYs 2016, 2015, 2014...

    - Same code, different input

# ACKNOWLEDGMENTS

**Division of Statistical Analysis and Reporting**

- Katrina Pearson

- Deepshikha RoyChowdhury

- Lindsey Scott

- Charles Wu

**Civic Digital Fellowship – NIH**

- Jessica Mazerik

- Amit Rajesh

**Funding Sources**



National Institutes of Health
*Turning Discovery Into Health*

The Washington Center

coding it forward >