# Building Scientific Publication Profiles for U.S.-trained Doctorate Recipients

Haoyi Wei[1], Eric Livingston[2], Christina Freyman[3], and Wan-Ying Chang[1]

[1]National Center for Science and Engineering Statistics, National Science Foundation

[2]Elsevier

[3]Evaluation and Assessment Capability Section, National Science Foundation

# Disclaimer

Working papers are intended to report **_exploratory_** results of research and analysis undertaken by the National Center for Science and Engineering Statistics at the National Science Foundation (NSF). Any opinions, findings, conclusions, or recommendations expressed in this work do not necessarily reflect the views of NSF. This work is being presented to inform interested parties of ongoing research or activities and to encourage further discussion of the topic.

# Brief Summary

**Research Goal:**

Identify the Scopus publication records for U.S.-trained doctorate recipients in science, engineering, and health fields.

**Data:**

Clarivate's Survey of Doctorate Recipients (SDR) - Web of Science (WoS) linkage

Elsevier's SDR – Scopus linkage.

**Method:**

Machine Learning

**Results:**

Finalized SDR-Scopus linkage.

**Contribution:**

This work is an important contribution to the development of a U.S. Science & Engineering Enterprise Data Network.

# Motivation

- Why do we build publication profiles?
    - To advance understanding of scientific research and the impact made by individual researchers.
    - To investigate the relationships between scientific productivity and various author attributes.
    - To access the nation's strategic investment in doctoral training to inform science policies.
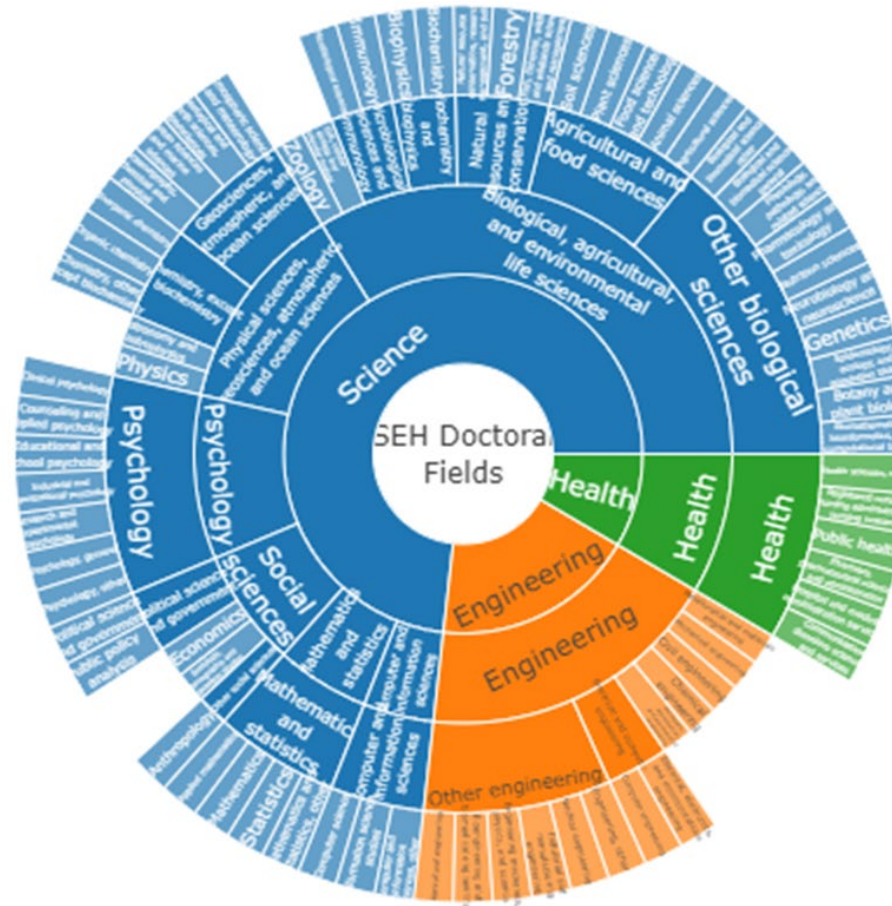


Citation Bias?                    Credit: Biomedical Odyssey

# Data

Surveys and bibliometric data

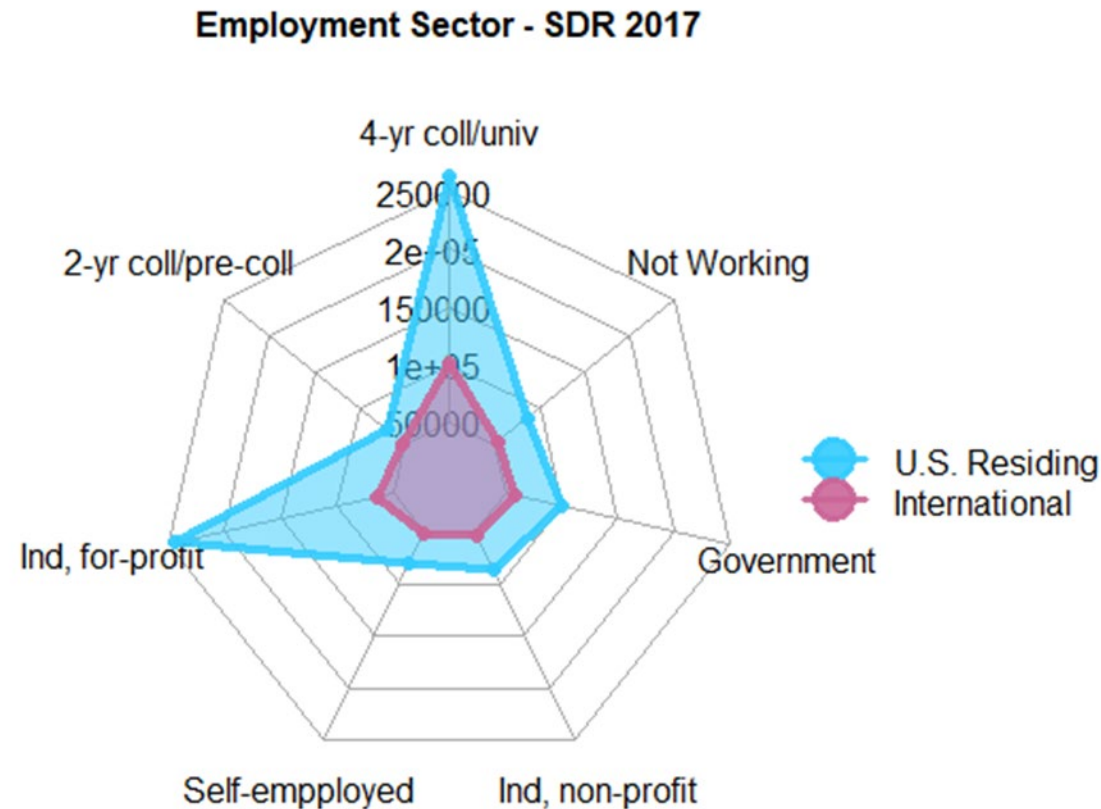# Surveys on U.S.-Trained Doctoral Researchers

Survey of Earned Doctorates (SED): annual census of research PhDs



SOURCE: National Center for Science and Engineering Statistics, Survey of Earned Doctorates and Survey of Doctorate Recipients, 2017.

# Surveys on U.S.-Trained Doctoral Researchers

Survey of Doctorate Recipients (SDR): biennial sample survey on science, engineering, and health doctorate degree holders



**Employment Sector - SDR 2017**

SOURCE: National Center for Science and Engineering Statistics, Survey of Earned Doctorates and Survey of Doctorate Recipients, 2017.

# Longitudinal data from SED to SDR

**Education**

1st bachelor, master, research PhD, up to 5 degrees (year, institution, place, field of study); dissertation field, financial support, debt

**Postgraduation Plans**

Country/State intend to live, taking a postdoc, employment status/commitment, employer type, salary, work activities

**Background**

Sex, marital status, dependents, parent's educational attainment, birthplace, citizenship, race/ethnicity, disability

**Employment Situation**

Labor force status, reasons for not working, year retired, principal employer, faculty rank, tenure status, principal job, work activity, salary, benefits, job satisfaction, federal support

**Past Employment**

**Other Work-Related Experiences**

**Recent Educational Experiences**

**Demographic Information**

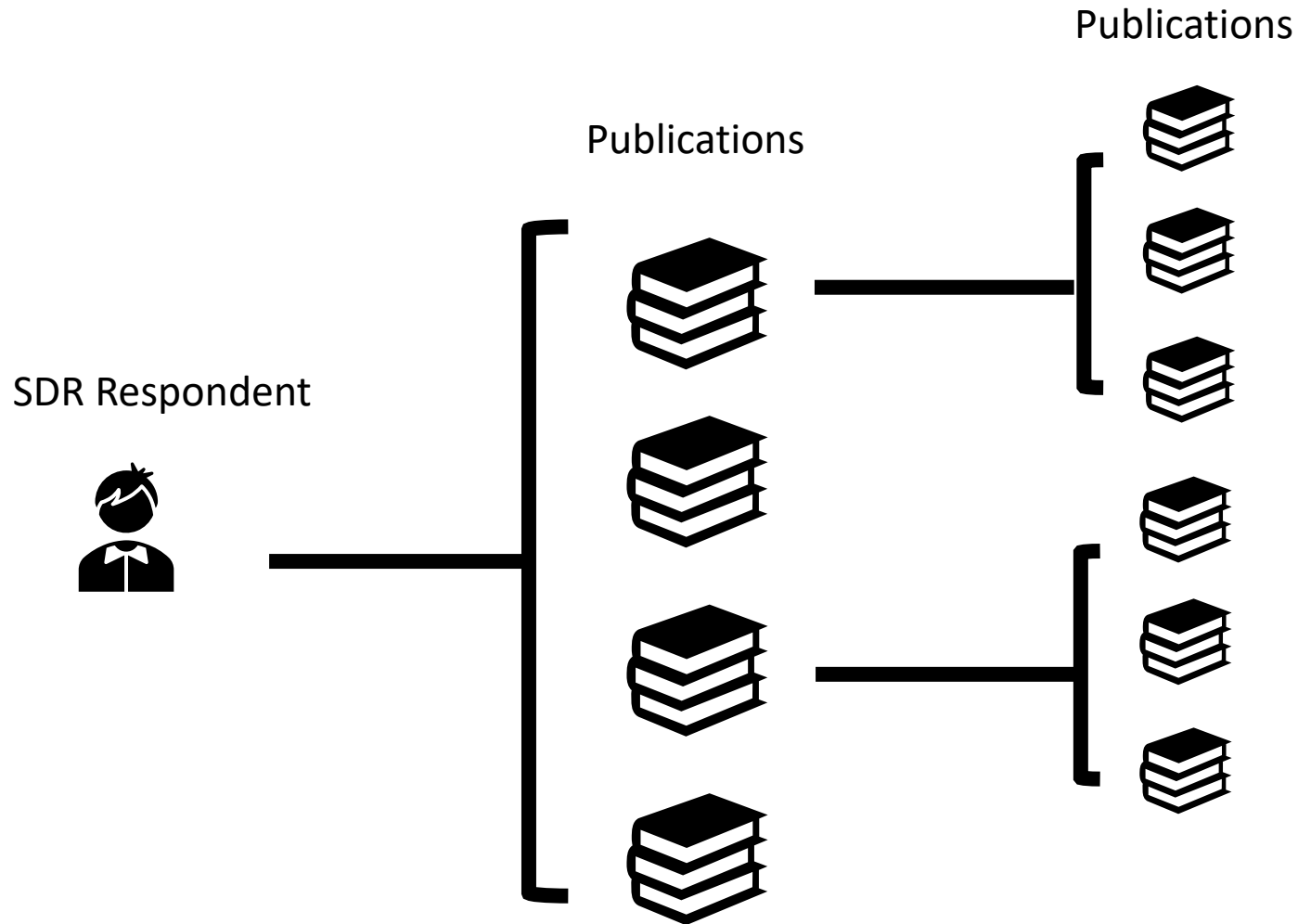Spouse working status, living with children, residing location, citizenship/visa type

# Publication Databases





The Web of Science (WoS) and Scopus are two leading databases providing reference and citation data from academic journals, conference proceedings, and other documents in various academic disciplines.
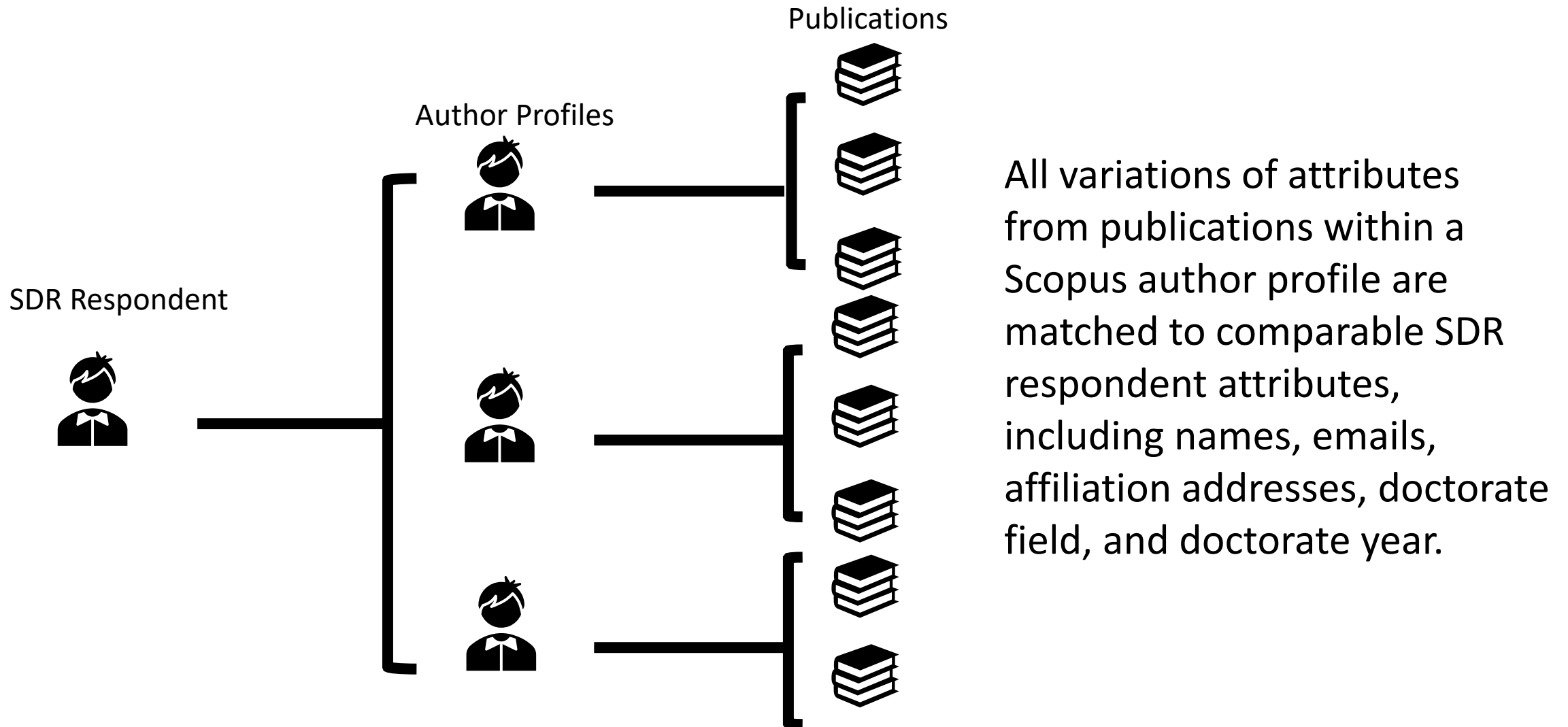
# Linkage 1: SDR-WoS



Two stages of machine learning approaches implemented sequentially to address
1. survey-to-publication linking
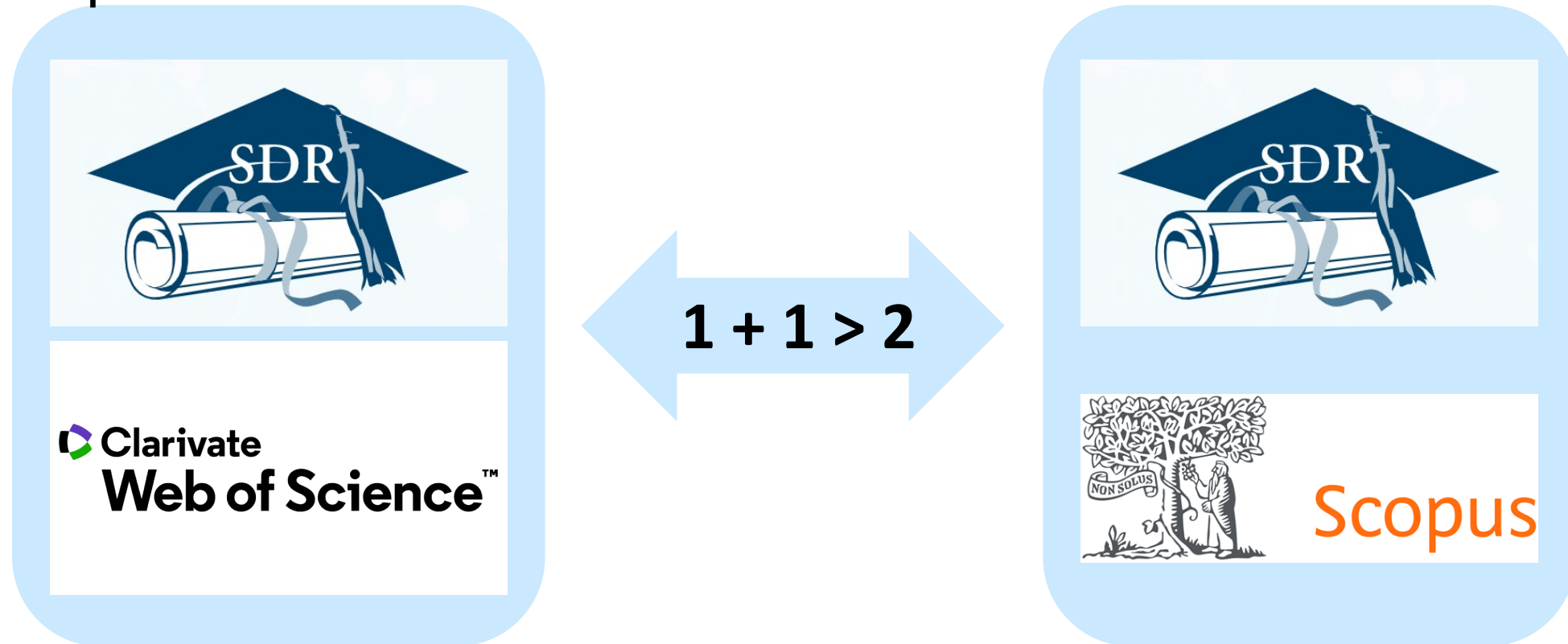2. Publication-to-publication linking

# Linkage 2: SDR-Scopus (Preliminary)



Publications

Author Profiles

SDR Respondent

All variations of attributes from publications within a Scopus author profile are matched to comparable SDR respondent attributes, including names, emails, affiliation addresses, doctorate field, and doctorate year.

# Methods

Connecting two independent data linkages

# Approach

- Connect two independent large scale data linkage to gain coverage, data quality, and cost-effectiveness for future expansion.
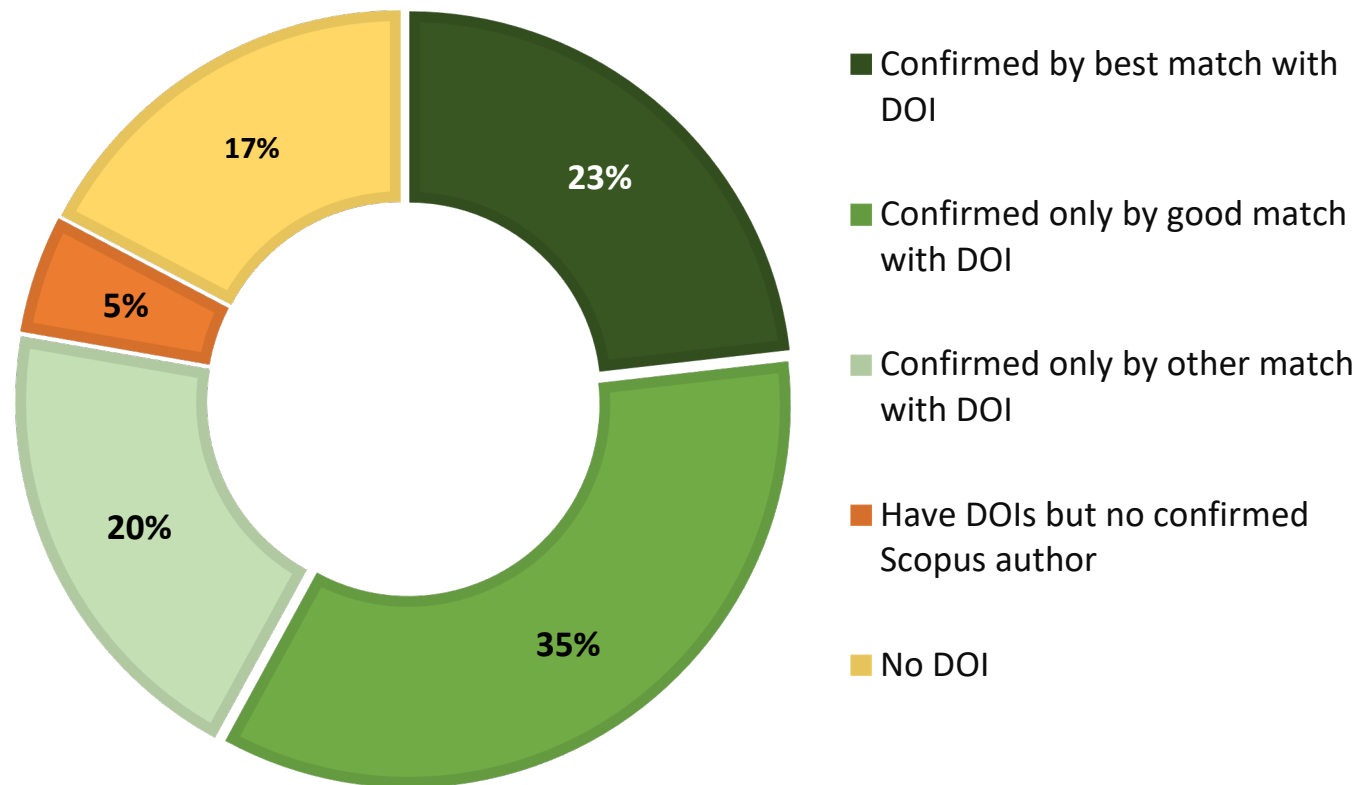


**1 + 1 > 2**

# Join SDR-WoS and SDR-Scopus

Digital Object Identifier (DOI) is the bridge to connect matched publications

**SDR-WoS Matched Pairs**

DOI=00000

REFID=123 — DOI=11235

No DOI

54% of publication have DOI

**Merge**

**SDR-Scopus Candidate Pairs**

Author ID=55 — DOI=11235 / DOI=11111

REFID=123

Author ID=66 — DOI=2222

Author ID=77 — DOI=3333 / DOI=4444

Author ID=88

DOI=5555

# DOI data can effectively confirm Scopus authors

## SDR-WoS authors



Legend:
- **Confirmed by best match with DOI** — 23%
- **Confirmed only by good match with DOI** — 35%
- **Confirmed only by other match with DOI** — 20%
- **Have DOIs but no confirmed Scopus author** — 5%
- **No DOI** — 17%

Although DOI is known for only half of the SDR-WoS publications, 83% of SDR-WoS authors have publication DOI data and 78% of SDR-WoS authors find DOI-confirmed Scopus authors in the candidate pool.

# Machine learning

Constructing training and evaluation data

# Construct training data – positive sample

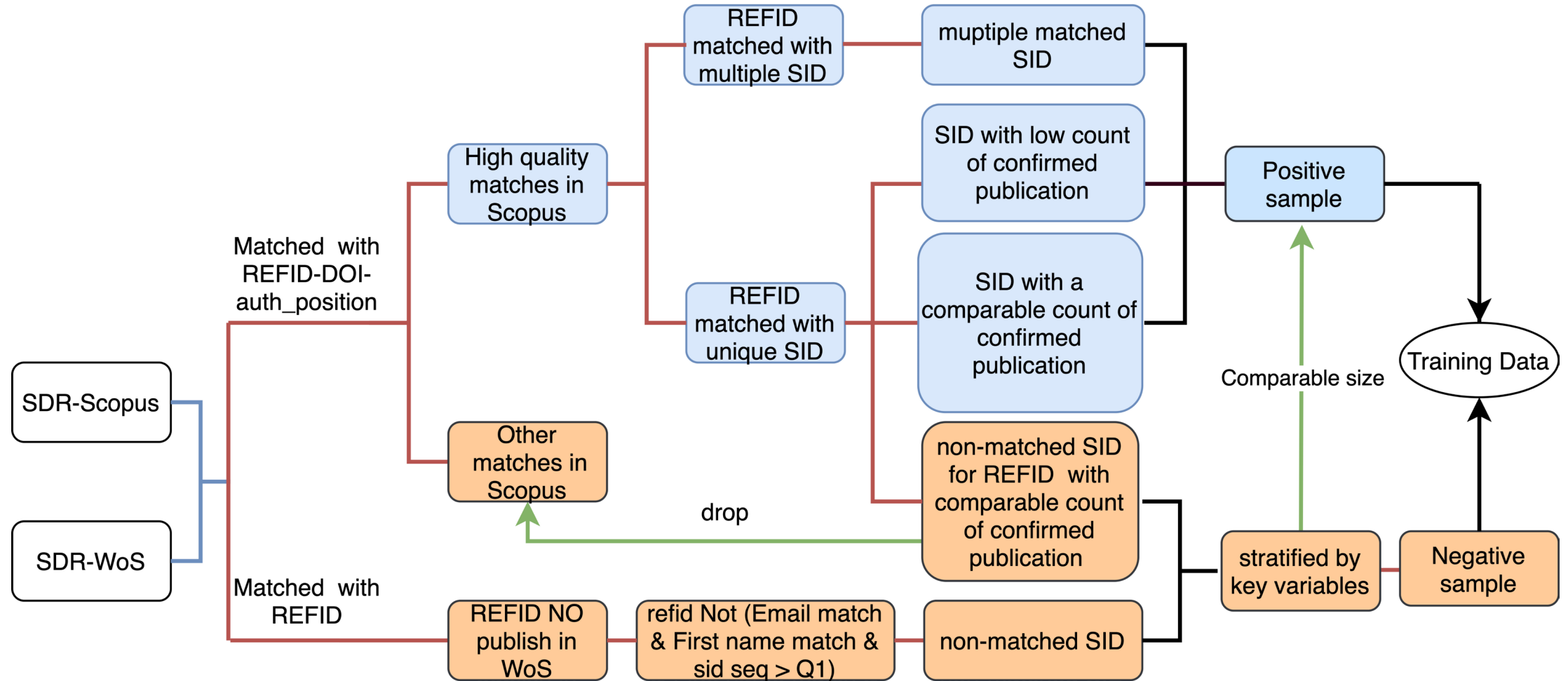- Keep REFID-author ID pairs with confirmed DOIs appeared in high-quality SDR-WoS matches.

# Construct training data – negative sample

-   Wrong matches within same refid: remaining (SDR, Scopus author ID) pairs under a respondent having a confirm pair with comparable total matched publications.
-   Non-authors: Candidate pairs under a respondent with no matched publications in WoS.

# Constructing the training data

# Machine Learning

## Training and evaluation

- Split the training data into training and test sets.

## Predictors

- Background and employment outcome: SED and SDR survey data
- Richness of source data for matching: quantity and quality of matching keys
- Scores of similarity: component scores of SOLR query

## ML methods tried

- Logistic regression
- Regression tree
- Random Forest

# Precision and Recall

Precision is the ratio between the correct predictions and the total predictions.

Recall is the ratio between the correct predictions and the total number of correct items in the set.

# Evaluate ML predictions



Preliminary results showing comparable performance among the three ML methods. Random Forest model performs slightly better in the prediction threshold of [0.5, 0.8]
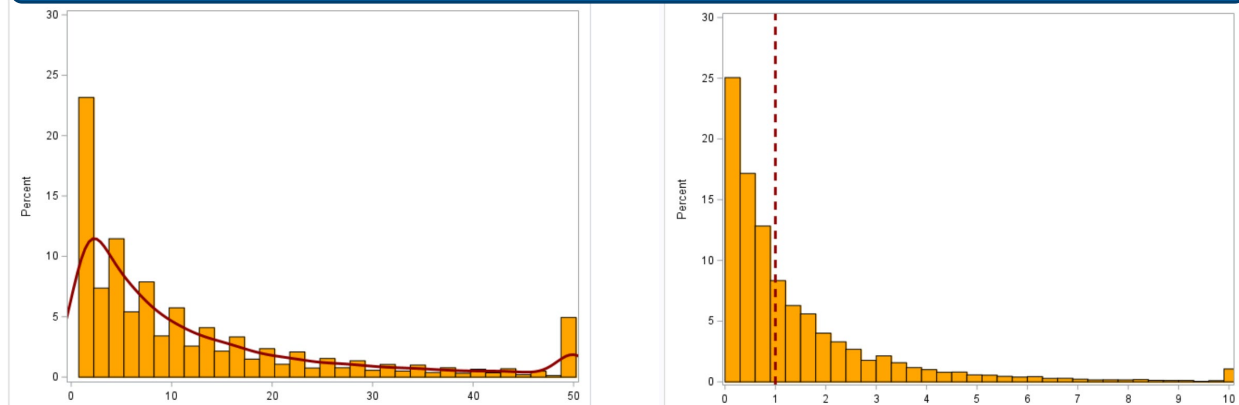
SOURCE: National Center for Science and Engineering Statistics, Survey of Doctorate Recipients linked to Web of Science Bibliometric Database and Scopus Author Profiles, 2021.

# Preliminary findings

ML predictions

# Match rates: WoS vs. Scopus



Overall match rate: WoS (72%), Scopus (65%)

SOURCE: National Center for Science and Engineering Statistics, Survey of Doctorate Recipients linked to Web of Science Bibliometric Database and Scopus Author Profiles, 2021.

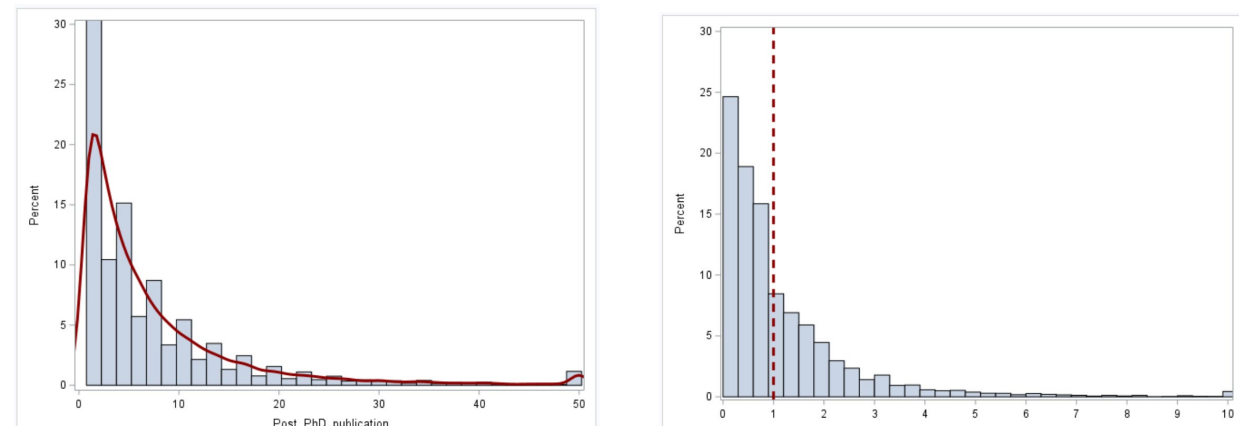# Post-graduation publications: 2008-2012 graduates



Preliminary Scopus matches show similar rate of post-PhD publications than SDR-WoS matches. Scopus author profiles contain more post-PhD publications

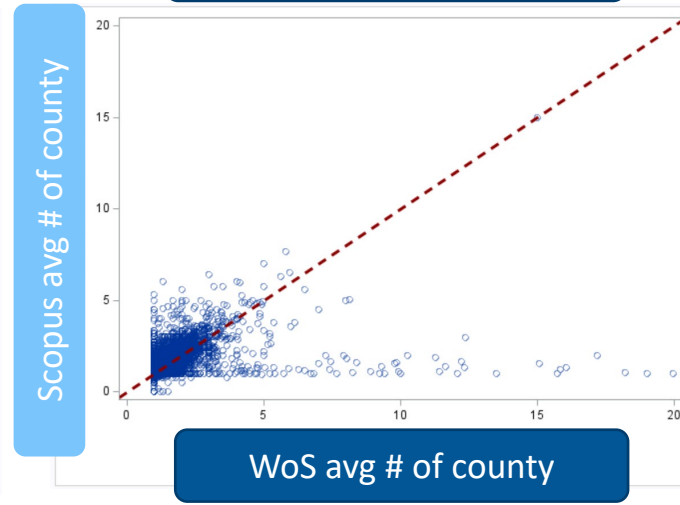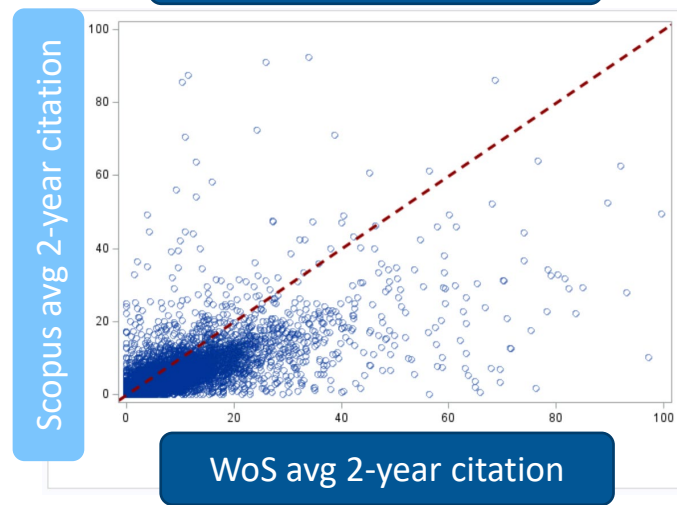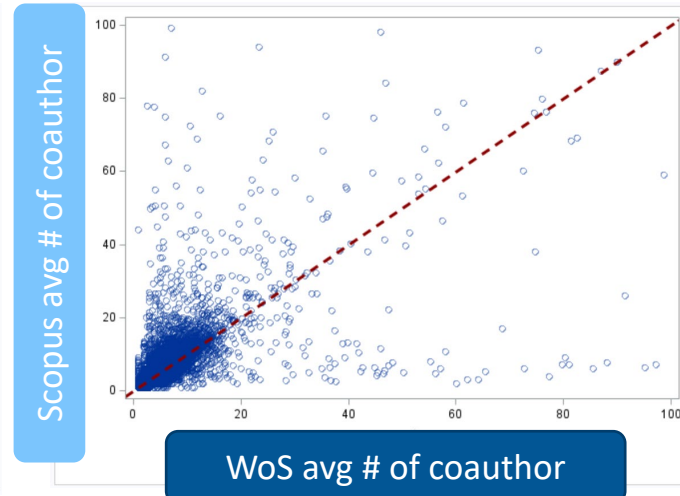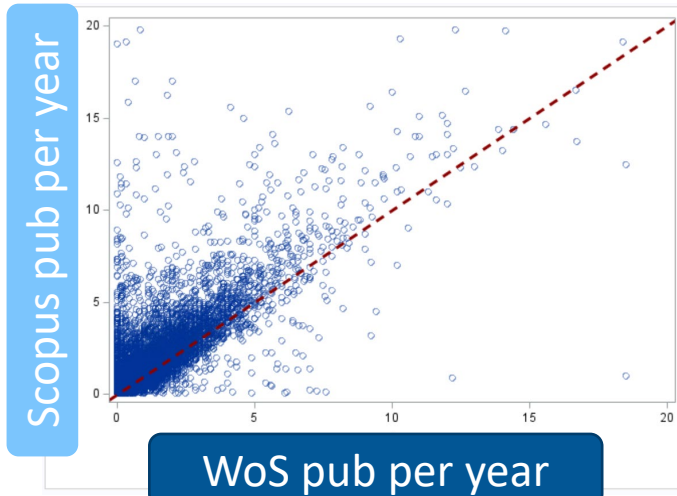# Productivity, collaboration, and impact



WoS pub per year



WoS avg # of coauthor



WoS avg 2-year citation



WoS avg # of county

At individual author level, more matched publication per year, more coauthors, more countries of affiliation on average are observed from Scopus matches. Average 2-year citations are higher from WoS matches.



# of pub in Scopus/# of pub in WoS
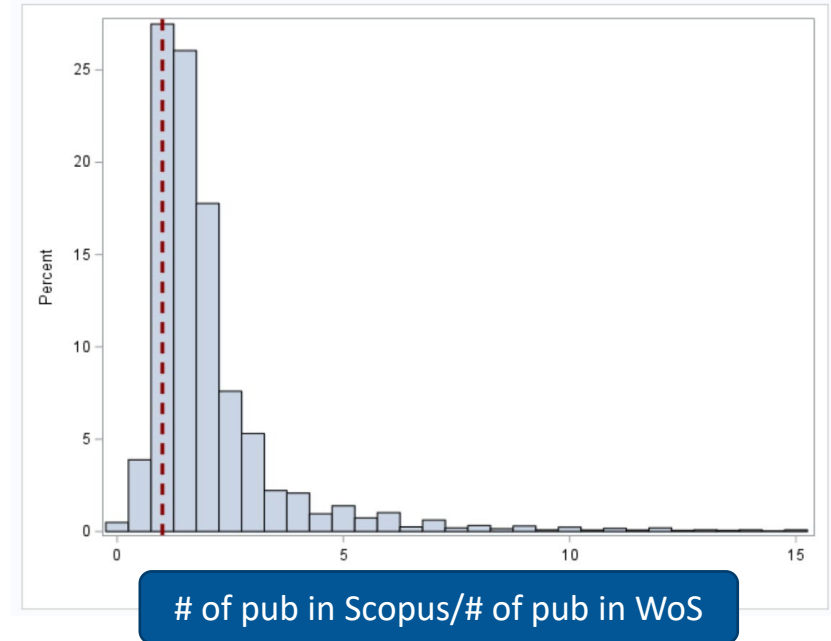
SOURCE: National Center for Science and Engineering Statistics, Survey of Doctorate Recipients linked to Web of Science Bibliometric Database and Scopus Author Profiles, 2021.

# Next steps

- ML modeling refinement

- Coverage bias analysis

- Matching quality evaluation

- Linkage updates

- Build use cases

Contact:   Haoyi Wei          howei@nsf.gov
           Wan-Ying Chang      wchang@nsf.gov