# Subject Indexing
# for Author Name Disambiguation

## Opportunities and Challenges

Cornelia Hedeler, Andreas Oskar Kempf, Jan Steinberg

# Outline

- Introduction
- Methods
- Subject indexing for author disambiguation
- Concrete use case at GESIS
- Experimental pre-study
- Multi-level approach:

  - Macro-level

  - Meso (intermediary/group)-level

  - Micro (individual)-level
- Conclusion
- Outlook

# Introduction (1/2)

- Increased demand for research monitoring
- Increased requirements on reporting systems of research activities
- Increased assessment of researchers/research institutions in terms of their impact
- Research information systems have become increasingly widespread – needed for institutional assessment procedures, accreditations, university rankings etc.

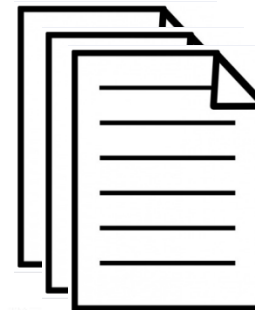> **>** Increased demand for quality-assured data on research activities

# Introduction (2/2)

**Which publication belongs to whom?**

Schmidt, Thomas[1]

Schmidt, Thomas[2]

Schmidt, Thomas[3]

Name string
*e.g. Schmidt, Thomas*

# Methods for author name disambiguation

- **Author grouping methods**
  - Using clustering techniques on a given data set on the basis of publication properties (co-authors, publication year etc.) to decide whether to group/subdivide publication records assigned to a certain author name.
- **Author assignment methods**
  - Directly assigning each publication to a given author by constructing a model that represents the author (e.g. the probabilities of an author publishing an article with other co-authors, in a given venue or using specific terms in the title of a publication).

# Publication properties

- Publication information:
  - Stream of the name
  - Co-author names (rate of co-authorship depends on the discipline)
  - Affiliations
  - Publication venue title
  - Publication year
  - Keywords from the title
- Additonal information:
  - Email addresses
  - Postal addresses
  - Data retrieved from the web

# Example of application (1/2)



> **Number/letter** = indication of co-authorship

# Example of application (2/2)

| Sonya Coleman | Sean S. Coleman   S. R. Coleman   S. S. Coleman   S |
| --- | --- |
| | John S. Coleman [T]   Nastaran S. Coleman   Robert S. |
| Bryan Scotney | Bryan A. Scotney [T] **1d 3a**   Bryan W. Scotney [ST] **1d 3a** |
| Bryan Gardiner | Bryan Gardiner [ST] **1d 2a 2b**   M. B. Gardiner   Paul H. B |

> **T** = indication of time proximity

> **S** = indication of topical closeness

# Concrete use case at GESIS



**GESIS Portal sowiport**

**Literature**
- 10 different data sources

- about 7,5 million publication records

**Future work:**
**>** Disambiguation of streams of author names
**>** Linking of individualized author name records to persistent identifiers

# Instruments for content cataloguing at GESIS

- Classification for the Social Sciences (CSS) > Disciplinary assignment

- Thesaurus for the Social Sciences (TSS) > Content information

gesis
Leibniz Institute for the Social Sciences

# Classification for the Social Sciences (CSS)

– **Field-/Disciplinary Classification in its present form since 1996 for GESIS databases of research literatur (SOLIS, SSOAR) and research projects (SOFIS)**

– **159 classes**

– **4 different hierarchical levels**

– **Variable amount of subclasses**

– **Indexing practice: 1 main class/notation and a variable amount of subordinate classes/notations**

# Classification cross-concordances

- Bilateral mapping between CCS and DDC (CSS > DDC / DDC > CSS )
(2012/2013)



| Start voc. | End voc. (main panels) | Rel. in total | Exact Match | Broader Match | Narrower Match | Null Rel. | Simple Matches | Multiple Matches (1:n) | Multiple Matches (n:1) |
|---|---|---|---|---|---|---|---|---|---|
| 159 | >27,000 | 169 | 45 | 15 | 89 | 1 | 20 | 41 | 16 |

# Thesaurus for the Social Sciences (since 1979)



- Translation into English, French, (Russian)

- About 8,000 subject headings

- About 4,000 non-descriptors/synonyms

- Classification scheme

- Cross-concordances to other vocabularies

- Indexing practice: 10-15 subject headings/document

# Thesaurus cross-concordances

- Major terminology mapping initiative (KoMoHe: 2004 – 2007) funded by the German Federal Ministry of Education and Research
- Mapped vocabularies i.a. AGROVOC, Medical Subject Headings (MeSH), Thesaurus for Economics (STW)

| Start voc. | End voc. | Rel. in total | Equiv. Rel. | BT-Rel. | NT-Rel. | Assoc. Rel. | Null-Rel. | Start terms | End terms | Term combi- nations |
|---|---|---|---|---|---|---|---|---|---|---|
| TSS | IAF (*SWD*) | 8,208 | 7,098 | 295 | 292 | 356 | 160 | 7,662 | 6,838 | 551 |
| IAF (*SWD*) | TSS | 9,432 | 6,276 | 1,831 | 134 | 640 | 594 | 8,890 | 5,556 | 182 |

TSS > IAF/IAF > TSS: Bilateral cross-concordance is mutually continuously developed

# Experimental pre-study

- **Research question:**
  Can topic information (subject headings, class notations) help to distinguish between different authors with the same name?

- **Data set:** Social Science Research Literature Information System (SOLIS) (since 1978):

  – Bibliographic records including content information on German social science literature (monographs, compilations, journal articles (300 journals) and grey literature)

  – 450,000 social science publications (Jan. 2014) from 1945-

  **>** At this point of our study we perceive name records in SOLIS as being fully disambiguated.

# Experimental pre-study – set up

- Approach: Longitudinal analysis of content information
    - Publication years: 1954 - 2013
    - „Profiling" of authors (more than two publications): 63,683 author names (81.14% of publications in single authorship)
    - Multi-level approach:
        - Macro-level:
        How discriminative/expressive are subject headings/class notations in general?

        - Meso- (group) level:
        How do topic distributions of research interests along a career look like?
            - Group 1 (5-10 years of publication activity): 16,108 author names
            - Group 2 (20-30 years of publication activity): 7,953 author names
            - Group 3 (40-50 years of publication activity): 482 author names

        - Micro-level:
        How do individual topic distributions of research interests along a career look like?
            - Group 1 (5-10 years of publication activity): one example
            - Group 2 (20-30 years of publication activity): two examples
            - Group 3 (40-50 years of publication activity): one example

# Subject headings

- Mean value of subject headings (whole thesaurus) for the average author: 48.11
- Mean value of authors per subject heading (whole thesaurus):  375.46
- for the selection of publications we have included in our pre-study frequencies of use (subject heading/author) range from 1 author (e.g. *official title* – ger. Amtsbezeichnung) to  29,599 authors (*historical development*)

**Preliminary results: Macro-level**

# Classification (1/2)

- Mean value of classes (**whole classification**) for the average author: 6.57

- Mean value of **core classes** (#1….) for the average author: 6.02

- Mean value of **aggregated classes** (first 3 digits of classification ID for whole classification) for the average author: 3.98

- Mean value of **aggregated core classes** (first 3 digits of classification ID for core classes) for the average author: 3.47

# Classification (2/2)

- Average number of author names per class (whole classification): 2,631
- Average number of author names per class (**core classes**): 2,655
- Average number of author names per **aggregated class** (first 3 digits of class-ID – whole classification): 7,908
- Average number of author names per **aggregated class** (first 3 digits of class-ID) **of core classes**: 10,421

> **>** The coverage of the core areas/classes is higher than for the more marginal areas/classes of the classification.

# Classification

- Average number of classes per author (**1st group: 5-10** years of publication):
  - Whole classification: 5.45
  - Aggregated classes: 3.61
- Average number of classes per author (**2nd group: 20-30** years of publication):
  - Whole classification: 10.98
  - Aggregated classes: 5.79
- Average number of classes per author (**3rd group: 40-50** years of publication):
  - Whole classification: 21.50
  - Aggregated classes: 9.31

# Subject headings

- Average number of different subject headings per author (**1st group: 5-10** years of publication):
  - 37.87
- Average number of different subject headings per author (**2nd group: 20-30** years of publication):
  - 84.48
- Average number of different subject headings per author (**3rd group: 40-50** years of publication):
  - 191.92

# Example 1st group: P. K. F.
# 10 years of publication activity)/29 publications

**Distribution of most frequently assigned classes**



Legend:
- **Philosophy of Science**
- **Sociology of Science**
- Philosophy & Religion
- Cultural Sociology
- Fundamentals of the Social Sciences
- General Concepts & Major Theories in the S.S.

**Preliminary results: Micro-level**

# Example 1st group: P. K. F.
# 10 years of publication activity)/29 publications

**Distribution of most frequently assigned subject headings**



**Science**
**Knowledge**
Theory
Logic
Experience
Philosophy of science

**Preliminary results: Micro-level**

# Example (a) 2<sup>nd</sup> group: S. K.
# 28 years of publication activity/17 publications

**Distribution of most frequently assigned classes**



**Population Studies**
**Demography**
**Research Design**
**Basic Research/General**
**Concepts of Demography**
Economics
Methods & Techniques of Data
Collection

**Preliminary results: Micro-level**

# Example (a) 2nd group: S. K.
# 28 years of publication activity/17 publications

**Distribution of most frequently assigned subject headings**



**Statistics**
*Germany*
Labor force
**Mortality**
Working hours
Gainful employment

**Preliminary results: Micro-level**

# Example (b) 2nd group: A. N.
# 24 years of publication activity/146 publications

**Distribution of most frequently assigned classes**



General Sociology
General Concepts & Major Theories
Macrosociology
Philosophy & Religion
Cultural Sociology
Political Sociology

# Example (b) 2nd group: A. N. 24 years of publication activity/146 publications

**Distribution of most frequently assigned subject headings**



**System theory**
**Society**
**Luhmann, N.**
**Communication**
**Sociology**
**Modernity**

**Preliminary results: Micro-level**

# Example 3rd group: K. H.
# 41 years of publication activity/241 publications

**Distribution of most frequently assigned classes**



Legend:
- **Youth Sociology**
- **Social Psychology**
- **Sociology of Education**
- **Medical Sociology**
- **Social Problems**
- **Education & Pedagogics**

**Preliminary results: Micro-level**

# Example 3ʳᵈ group: K. H.
# 41 years of publication activity/241 publications

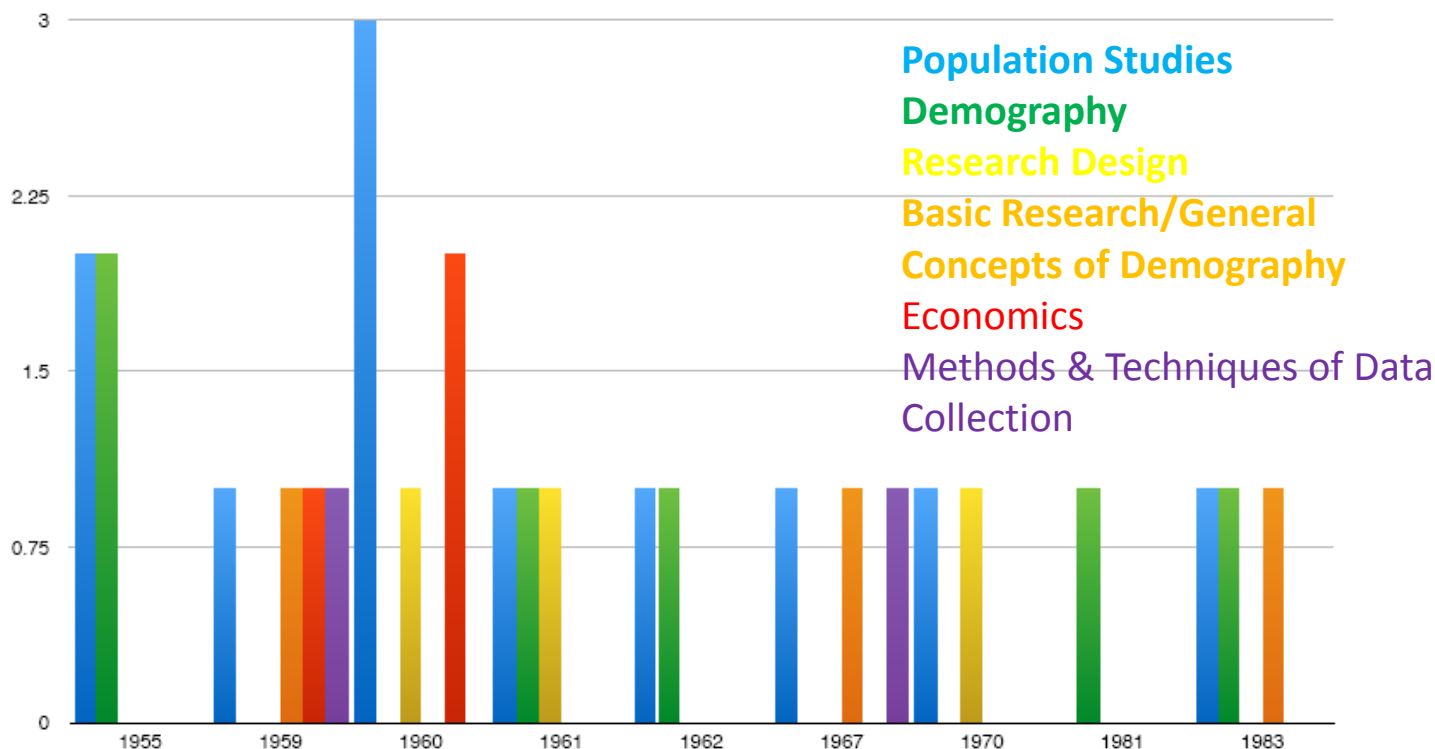**Distribution of most frequently assigned subject headings**



Legend:
- **Adolescent** (blue)
- *Germany* (green)
- **School** (yellow)
- **Child** (orange)
- **Health** (red)
- **Socialization** (purple)

# Preliminary results –
## for the German-speaking social sciences

- **Research question:**

  Can topic information (subject headings, classification) help to distinguish between different authors with the same name?
  OR: Is there subject continuity along a researcher´s career?

- **Macro-level of study:**

  – Speaking of an „average author" could only serve as an approximation to answer the research question.

  – The frequency of use differs enormously between different classes/subject headings.

- **Meso-level of study:**

  – The average number of classes/subject headings grows continuously along a researcher´s career. Using subject information for author name disambiguation therefore might appear rather fruitless. However,…

- **Micro-level of study:**

  – Considering the most frequently used classes/subject headings assigned to publications of an author could help to profile/ disambiguate an author name. Subject continuity could be identified. This already applies to authors with a rather short period of publication activity.

# Outlook

- **Future workflow**
  - ↓ 1st step of disambiguation: Preprocessing via standard algorithms on the basis of **database-internal** publication/author-centred reference information (publication properties like co-authors, year of publication etc.)
  - ↓ 2nd step of disambiguation: Consideration of classification and subject heading information for cases of doubt
  - ↓Linking person records to authority files: Use of **database-external** person-centred reference information (individualized/differentiated person records of the Integrated Authority File (IAF)) via mapping of bibliographic records and individualized person IDs on the basis of an overlap of subject information using cross-concordances between classifications (DDC/CSS) and thesauri (IAF/TSS).
    - Preliminary results: Differentiated person records in the IAF only cover a small section of author strings included in the GESIS portal sowiport (mainly German-speaking social scientists)
    - Future work: Taking cross-concordances between thesauri (TSS – IAF) and classifications (CSS – DDC) for author name disambiguation into account.

# Thank you very much for your attention.

**Contact**

Dr. Andreas Oskar Kempf (GESIS)
Andreas.Kempf@gesis.org

Dr. Cornelia Hedeler (University of Manchester)
chedeler@cs.manchester.ac.uk

Jan Steinberg (GESIS)
Jan.Steinberg@gesis.org

# References

- Ferreira, Anderson A.; Goncalves, Marcos André; Laender, Albert H. F. (2012) A Brief Survey of Automatic Methods for Author Name Disambiguation, *SIGMOD Record*, vol. 41, No. 2.
- Liu, Wanli; Dogan, Rezarta Islamaj; Kim, Sun; Corneau, Donald C.; Kim, Won; Yeganova, Lena; Lu, Zhiyong; Wilbur, John W. (2013) Author Name Disambiguation for PubMed, *JASIST*.
- Mayr, Philipp; Petras, Vivien (2008b): *Cross-concordances: terminology mapping and its effectiveness for information retrieval.* In: 74th IFLA World Library and Information Congress. Québec, Canada URL: www.ifla.org/IV/ifla74/papers/129-Mayr_Petras-en.pdf.
- Riege, Udo (1998) *Elektronische Version des Thesaurus und der Klassifikation Sozialwissenschaften*, IZ-Arbeitsbericht Nr. 14, GESIS IZ InformationsZentrum Sozialwissenschaften.
- http://www.gesis.org/unser-angebot/recherchieren/thesauri-und-klassifikationen/thesaurus-sozialwissenschaften/#c25531 (last accessed  July 1st, 2014).
- http://www.gesis.org/fileadmin/upload/dienstleistung/tools_standards/KassifikationSozialwissenschaften_Stand_Juli_2013_dt_en__2_.pdf (last accessed July 1st, 2014).