

Linked Open Research Data for Social Science

A concept registry for granular data documentation

**Jana Nebelin³, Pascal Siegers¹, Dagmar Kern¹, Antonia May¹, Fakhri Momeni¹, Andreas Daniel²,
Ben Zapilko¹, Claudia Saalbach³, Knut Wenzig³, & Jan Goebel³**

¹GESIS Leibniz-Institute for the Social Sciences, ²German Centre for Higher Education Research and Science Studies, ³German Socio Economic Panel (SOEP) at the German Institute for Economic Research

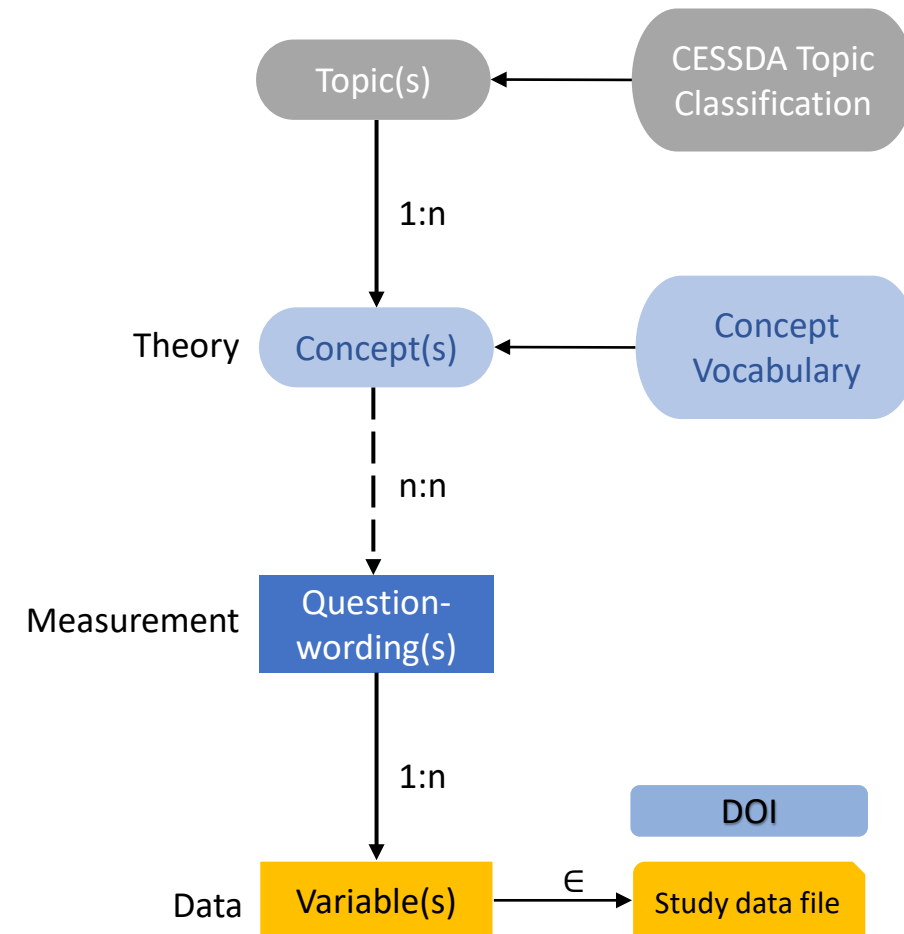
ESRA: July 2023 (Italy)

Outline of the presentation

- I. Introduction: the missing „link“ in data documentation(s)
- II. Concepts in Social Science Research
- III. Conceptualizing a Concept Registry
- IV. Lessons learned from the pilot Study
- V. Outlook: Establishing the LORD „pipeline“

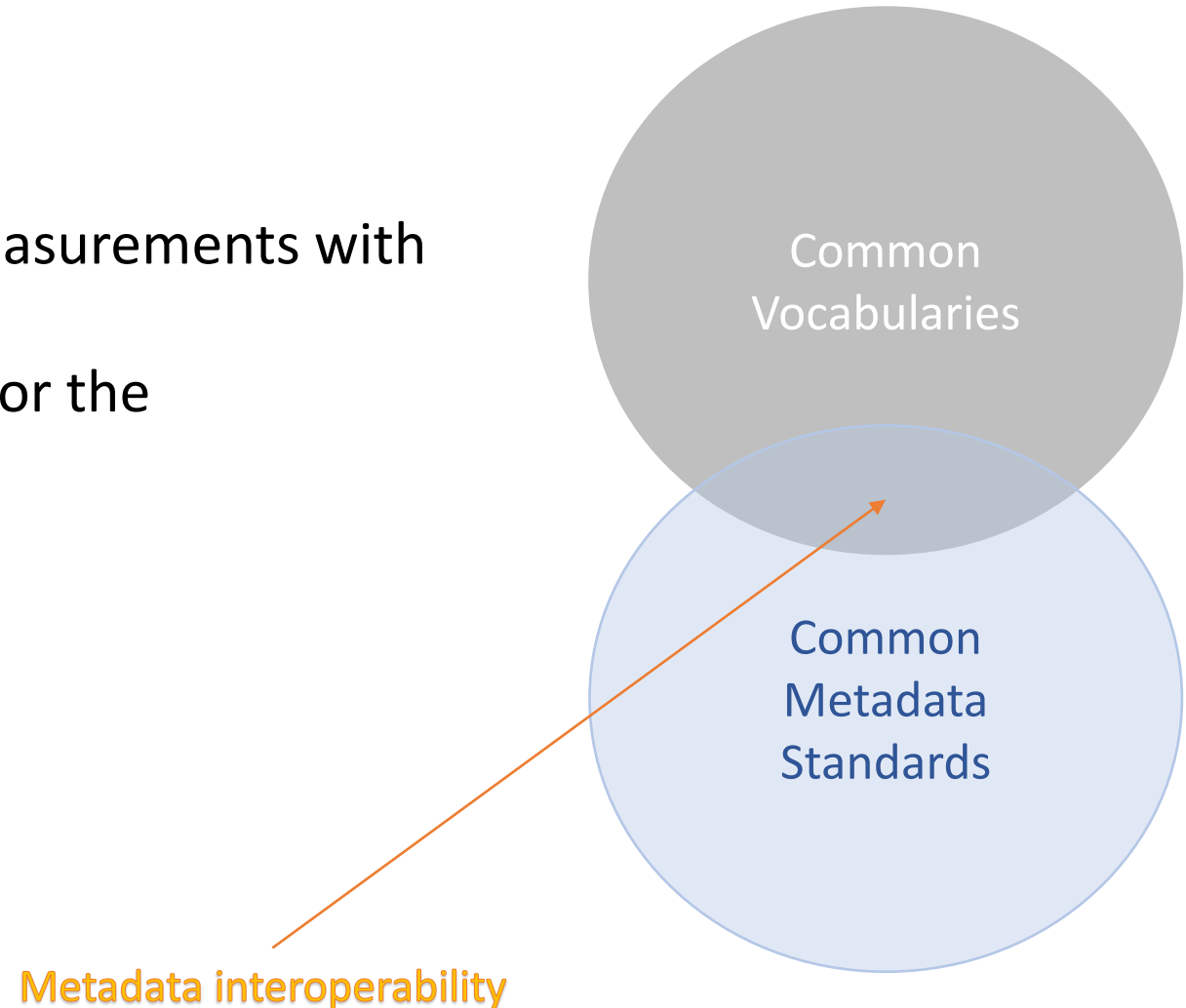
I. The missing „link“ in data documentation I

- General vocabularies for topics (e.g. ELSST, CESSDA Topic Classification)
- Extensive documentation of questions wordings (→ DDI)
- Extensive documentation of variables in data sets (labels, codes, code labels)
- Missing: often **no** information on theoretical concepts intended to measure
 - No concept vocabulary available for data documentation



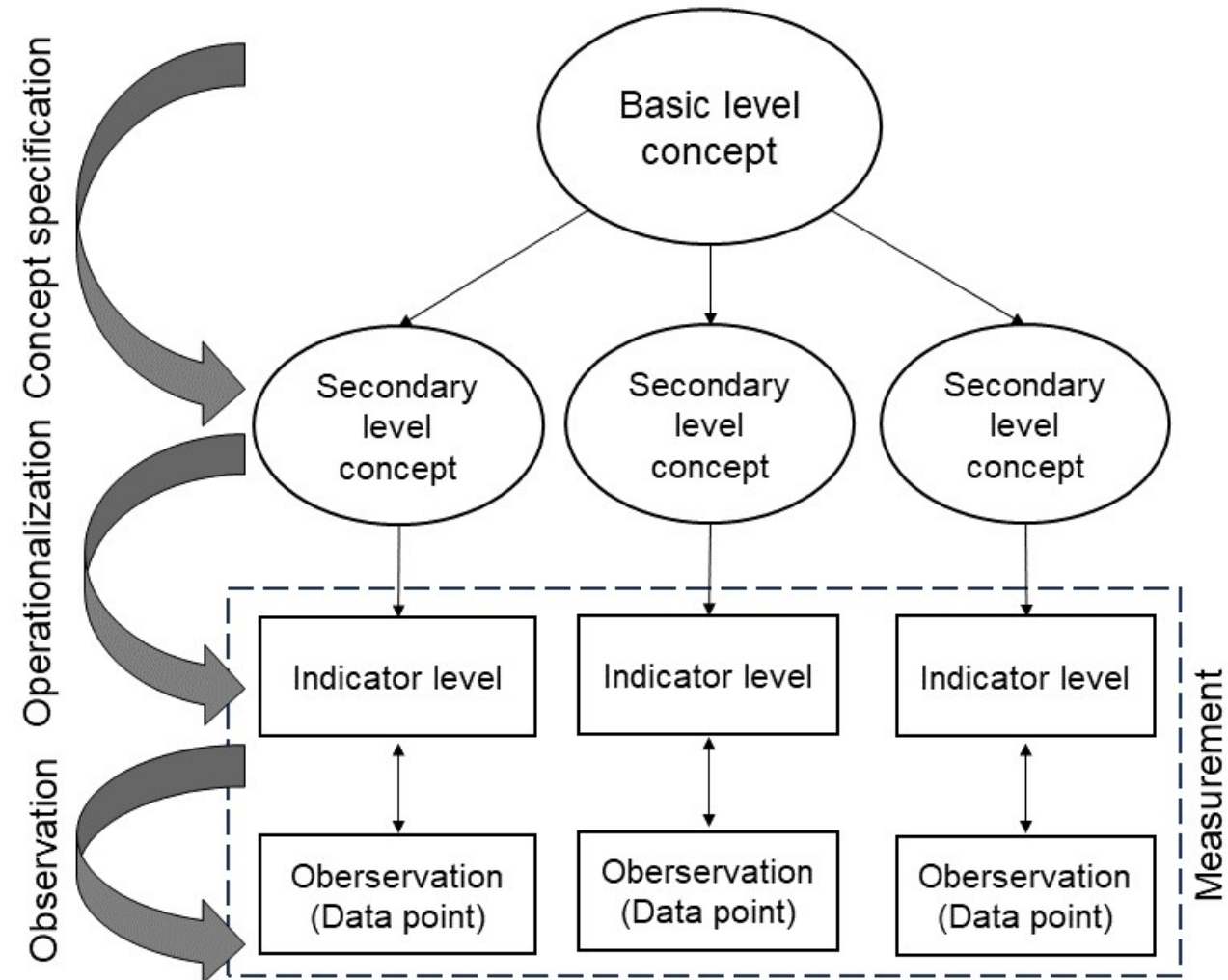
I. The missing „link“ in data documentation II

- Why concepts in documentation?
 - Supporting data search by linking measurements with concepts
 - Identifying different measurements for the same/similar concepts
 - > FAIRification of research data

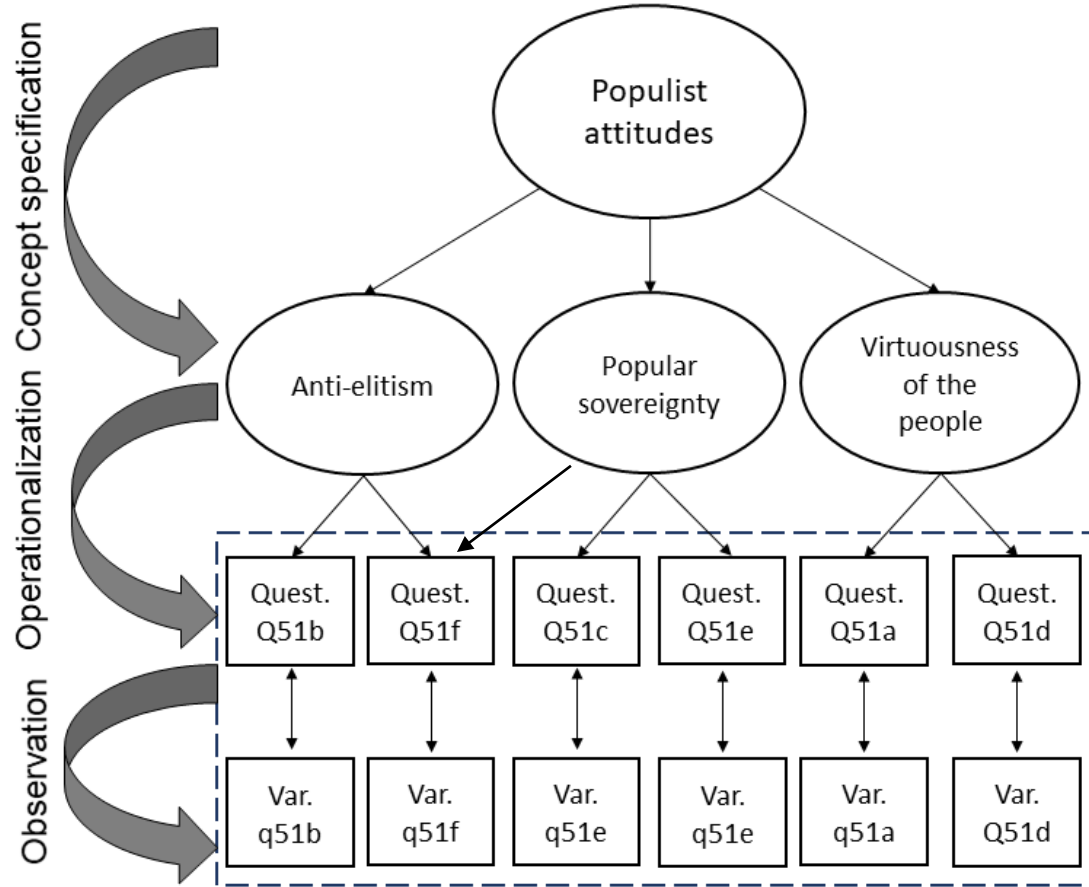


II. Concepts in Social Science research I

- Concepts are central elements of scientific language and knowledge representation
- Goertz (2006) distinguishes three levels of social science concepts
 - I. Basic level:** terminology used in theoretical propositions about reality
 - II. Secondary level:** Components of basic level concepts (dimensions)
 - III. Indicator level:** specifications for measurement



II. Concepts in Social Science research II



Example of indicator question (Q51d):

Question text:

Please say how much you agree or disagree with each of these statements.

Question Item

“Differences between the elite and the people are larger than the differences among the people.”

Answer scale

(1) Strongly agree; (2) Agree; (3) Neither agree nor disagree; (4) Disagree; (5) Strongly disagree.

Different measures for populism used in research practice

Diversity is the rule rather than the exception

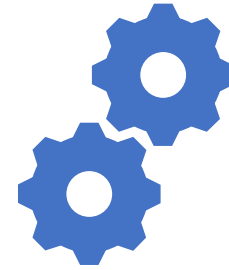
n:n relationships between concepts and measurements

III. Conceptualizing a Concept Registry



Construction principles:

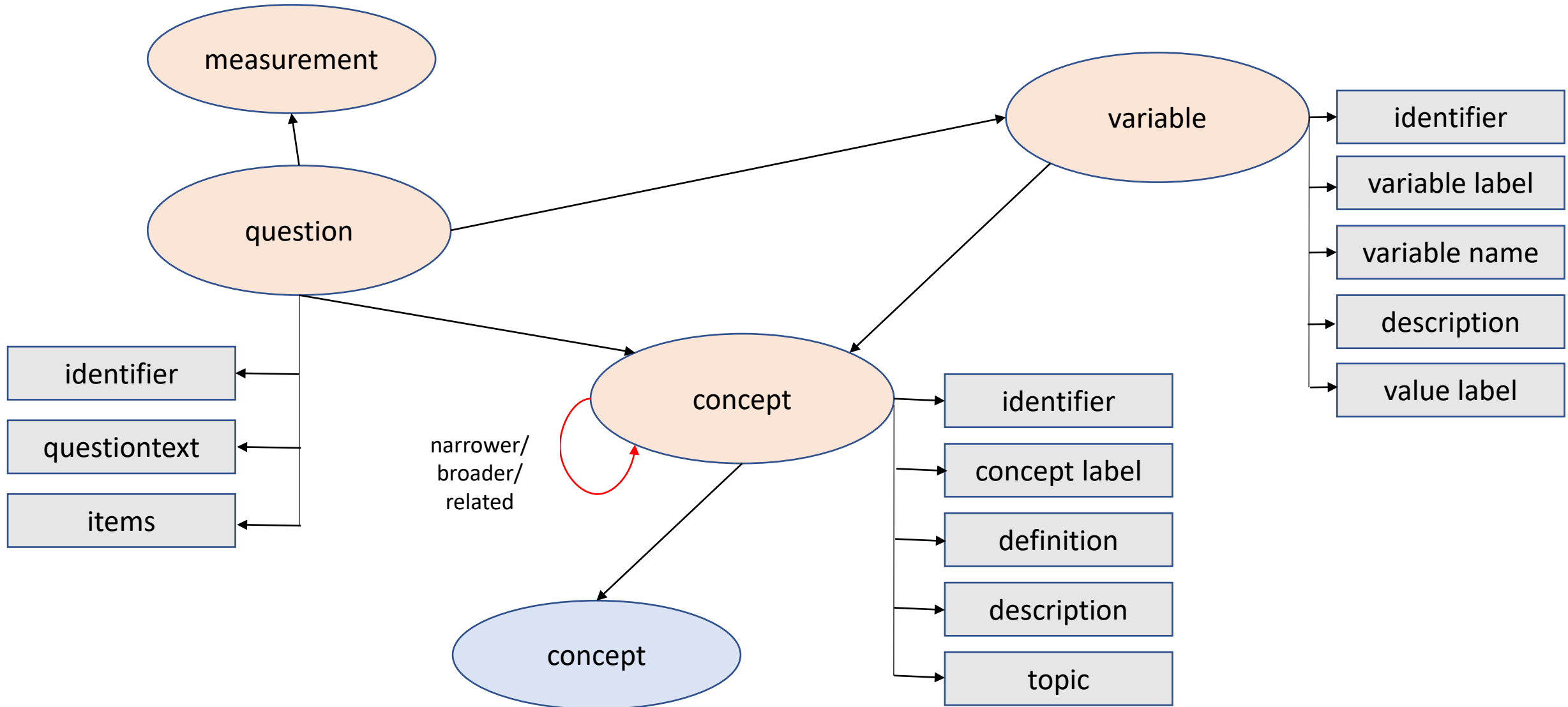
- I. Open and user driven development of the vocabulary
- II. Theory language
- III. Linking to existing vocabularies and LOD-Resources
- IV. Open interface(s) for re-use



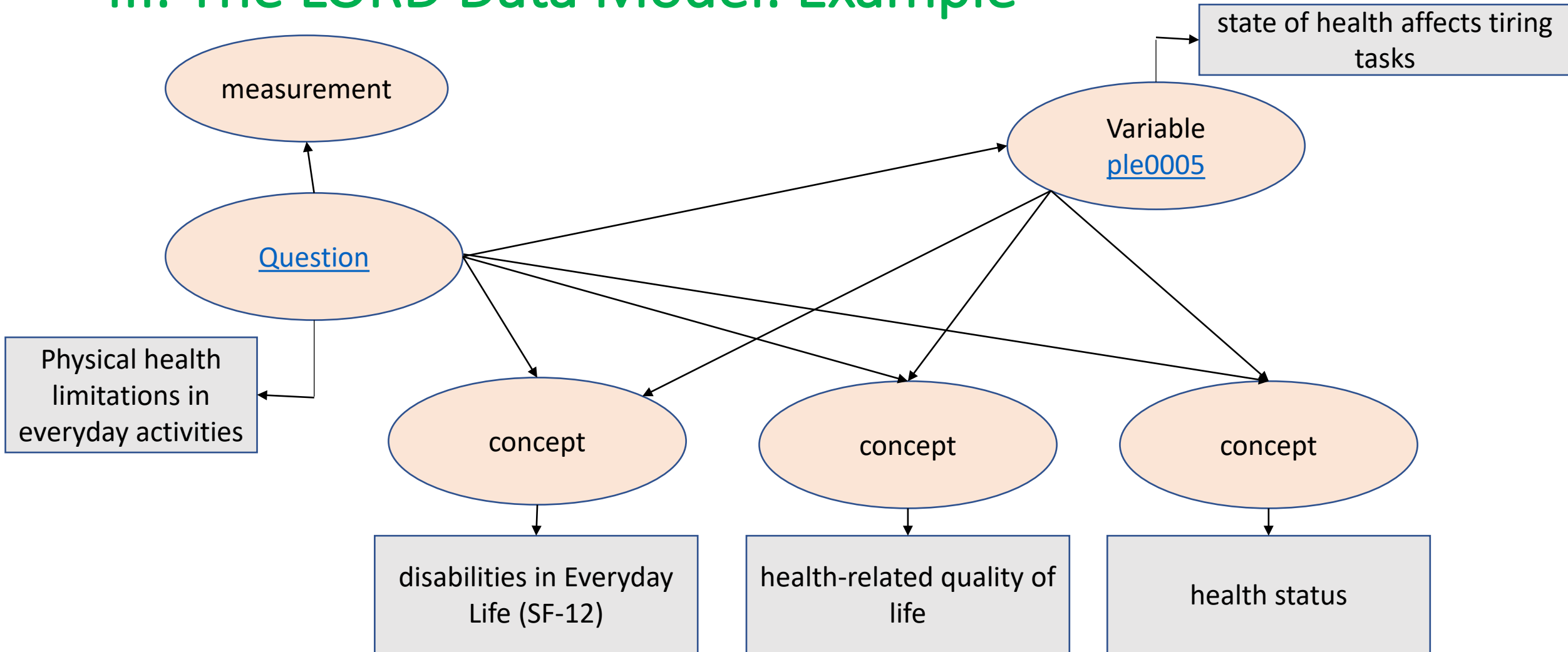
Components:

- I. Data model for the concept registry
- II. Annotation Tool (linking concepts to the measurement metadata)
- III. Triple Store

III. The LORD Data Model



III. The LORD Data Model: Example



III. The LORD Annotation Tool

- Displays question and variable metadata
- Allows to select a concept/topic from Thesaurus Social Sciences (TheSoz)
- New concepts are linked to the measurement and added to the concept registry
 - Automatic query extension supports selection of existing concepts
 - Several concepts can be linked to the metadata

The screenshot displays the LORD Annotation Tool interface, which is organized into several sections:

- Question ID:** ZA5274_Q348_Quelle
- Variable ID:** ZA5274_Varep03
- Fragetext:** Und Ihre eigene wirtschaftliche Lage heute?
- Frageitem:** (Empty text box)
- Antwortkategorien:** A dropdown menu showing a list of response categories: 1990: keine Teilnahme an Split 2 (Code 1 in spl90), Keine Angabe, Weiß nicht, Nicht erhoben 1980, 1988, Sehr gut, Gut, Teils gut / teils schlecht, Schlecht, Sehr schlecht.
- Variable label:** WIRTSCHAFTSLAGE, BEFR. HEUTE
- Wertelabel:** 1990: keine Teilnahme an Split 2 (Code 1 in spl90), Keine Angabe, Weiß nicht, Nicht erhoben 1980, 1988, Sehr gut, Gut, Teils gut / teils schlecht, Schlecht, Sehr schlecht.
- Konzept aus dem TheSoz vergeben:** A search bar with the placeholder text "search for a theme".
- Ausgewählte Konzepte:** (Empty list)
- Freie Konzepte, die nicht einem TheSoz Begriff zugeordnet werden können (optional):** A search bar with the placeholder text "Wahrn|", a list of suggestions including "Wahrnehmung der aktuellen eigenen...", "Wahrnehmung der aktuellen wirtscha...", and "Wahrnehmung der allgemeinen wirts...", and a button labeled "Freies Konzept hinzufügen".
- Kommen:** (Empty text box)

IV. Lessons Learned from the pilot study I

- Test annotation
 - German Socio-economic panel (GSOEP), German National Academics Panel Study (nacaps), and German General Social Survey (GGSS)
 - Each project partner annotated selection of questions from the three surveys (topics: health, income, migration etc.)
- Core questions for test:
 - Do annotations „overlap“?
 - Is there a *between concepts structure* „emerging“ from the annotations?

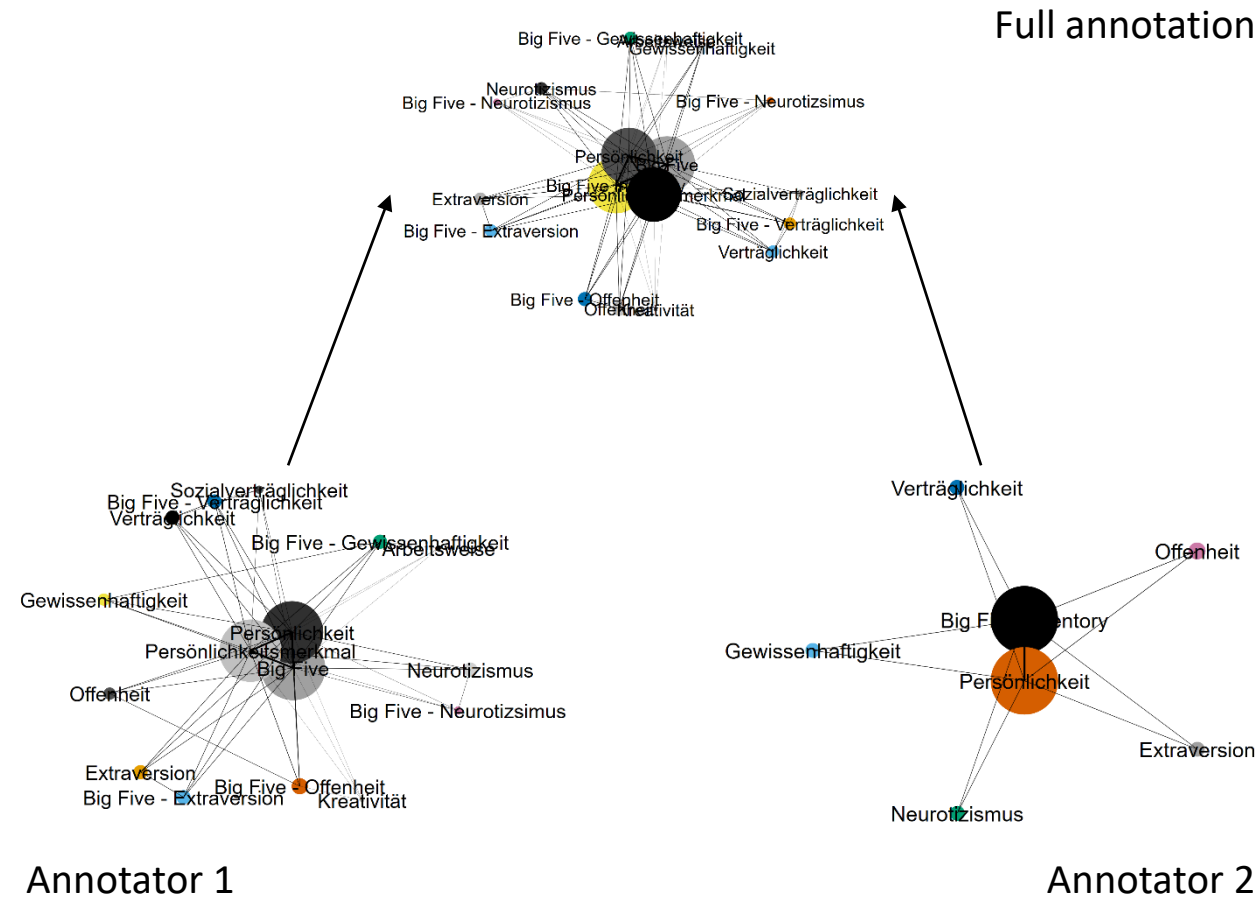
SOEP

ALBUS

nacaps 
National Academics Panel Study

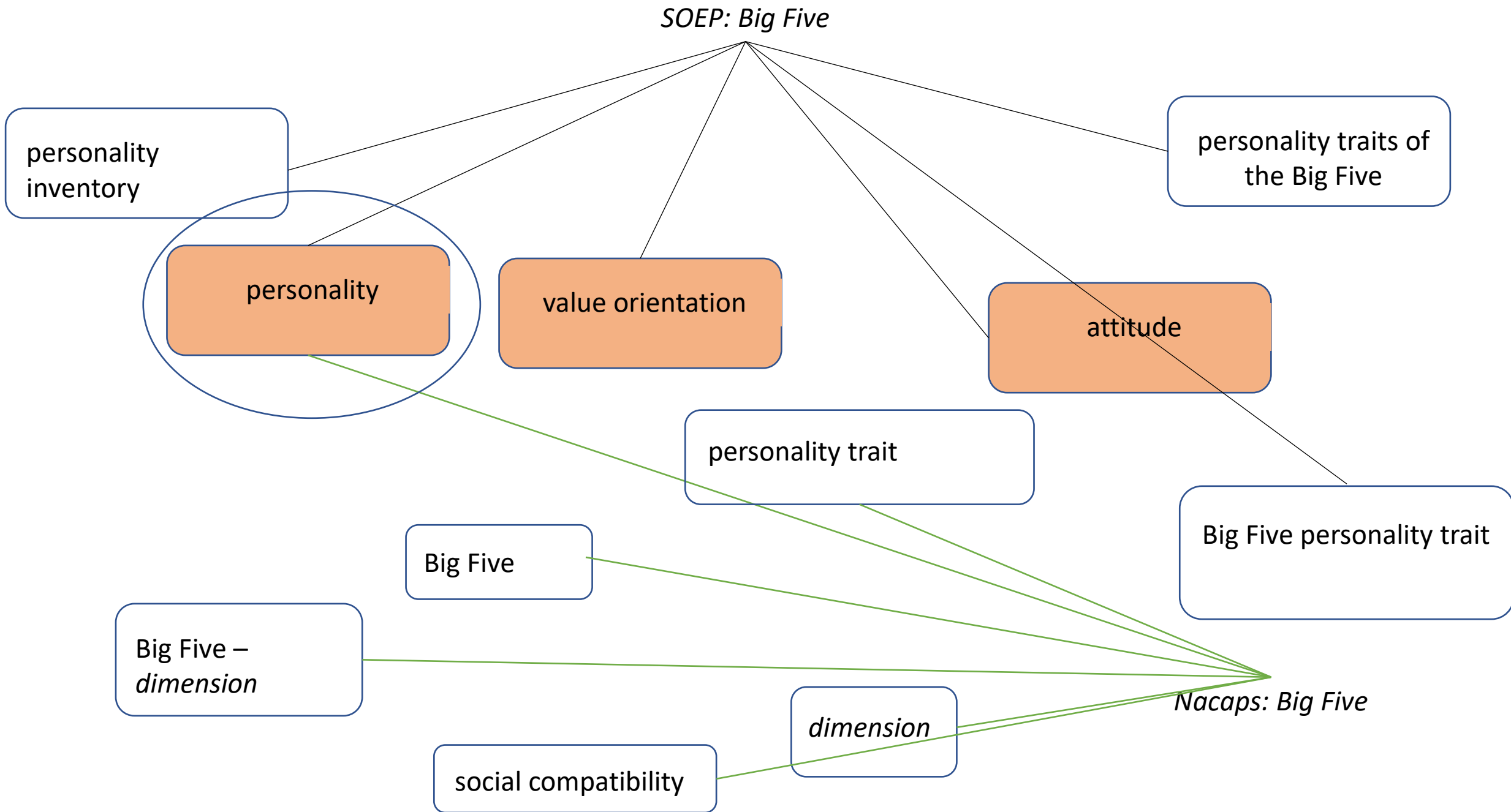
IV. Lessons Learned from the pilot study II

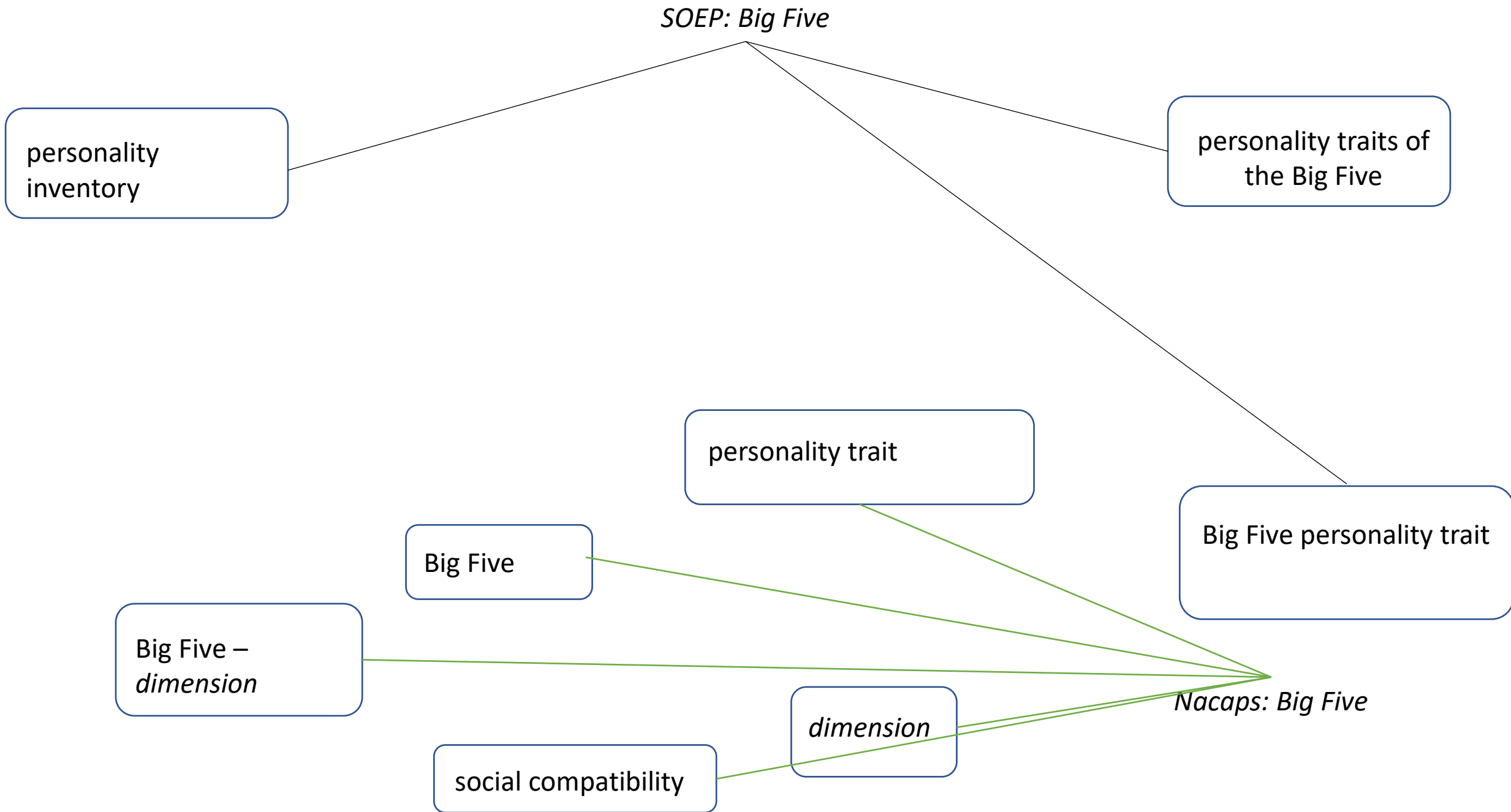
- Great diversity in individual annotation styles
- Results in large amount of different concept terms that cover very similar measurements
 - Non-substantive differences in concepts



IV. Lessons Learned from the pilot study II

- More general terms (topics) seemingly create links between concepts
- The current structure lacks the possibility to create links between concepts
 - Although this is part of the data model





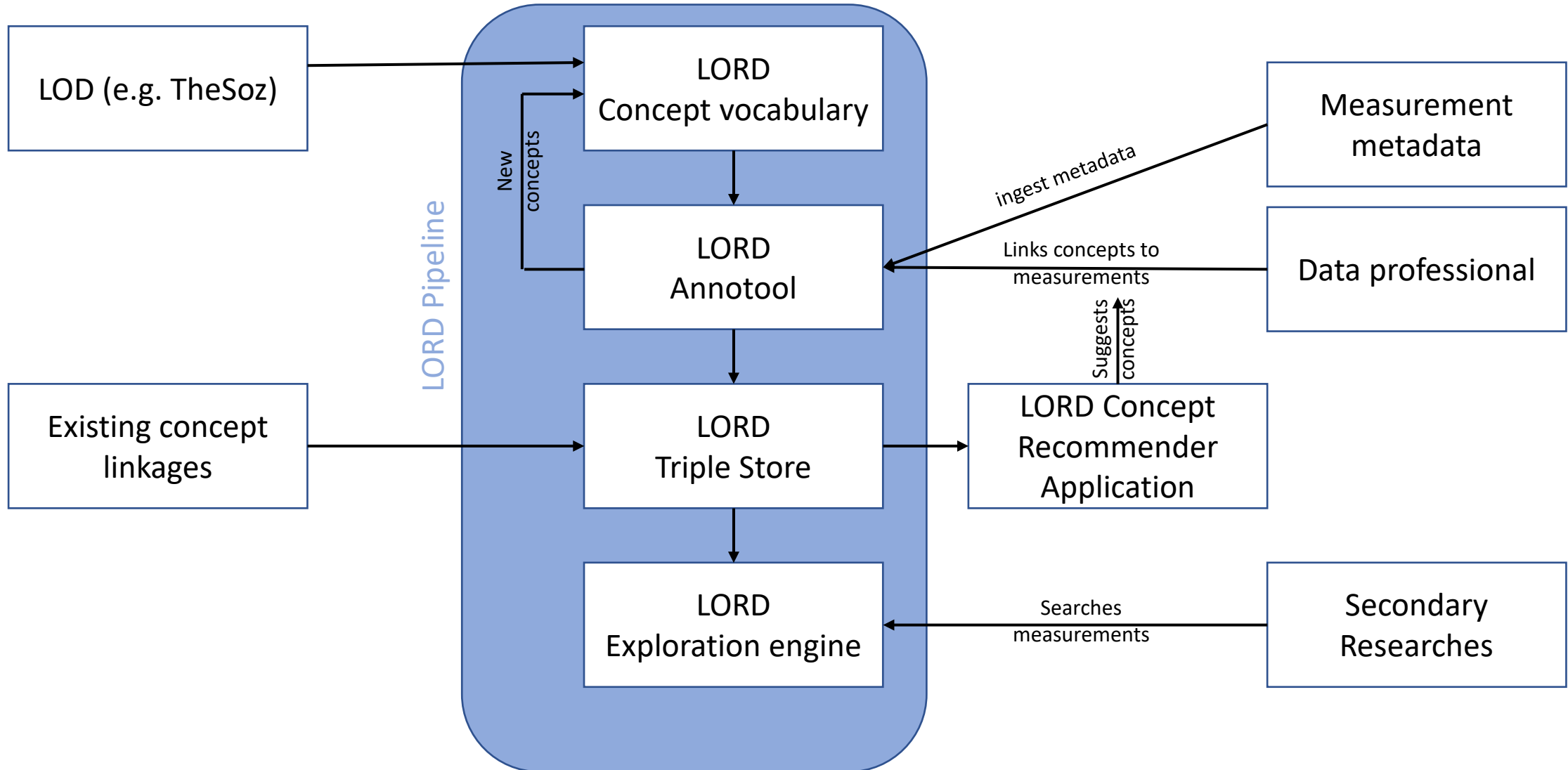
V. Outlook: developing the LORD „pipeline“

Current project only covers the exploration phase.

A user driven concept registry will require additional functionalities:

- Performant recommendation systems for concepts based on measurement metadata
 - Start phase: curated corpus of terms and relationships for the concept registry to improve recommender systems
- possibility to create links between concepts (not part of the current tools)
 - Graph based concept exploration engine

V. Outlook: developing the LORD „pipeline“



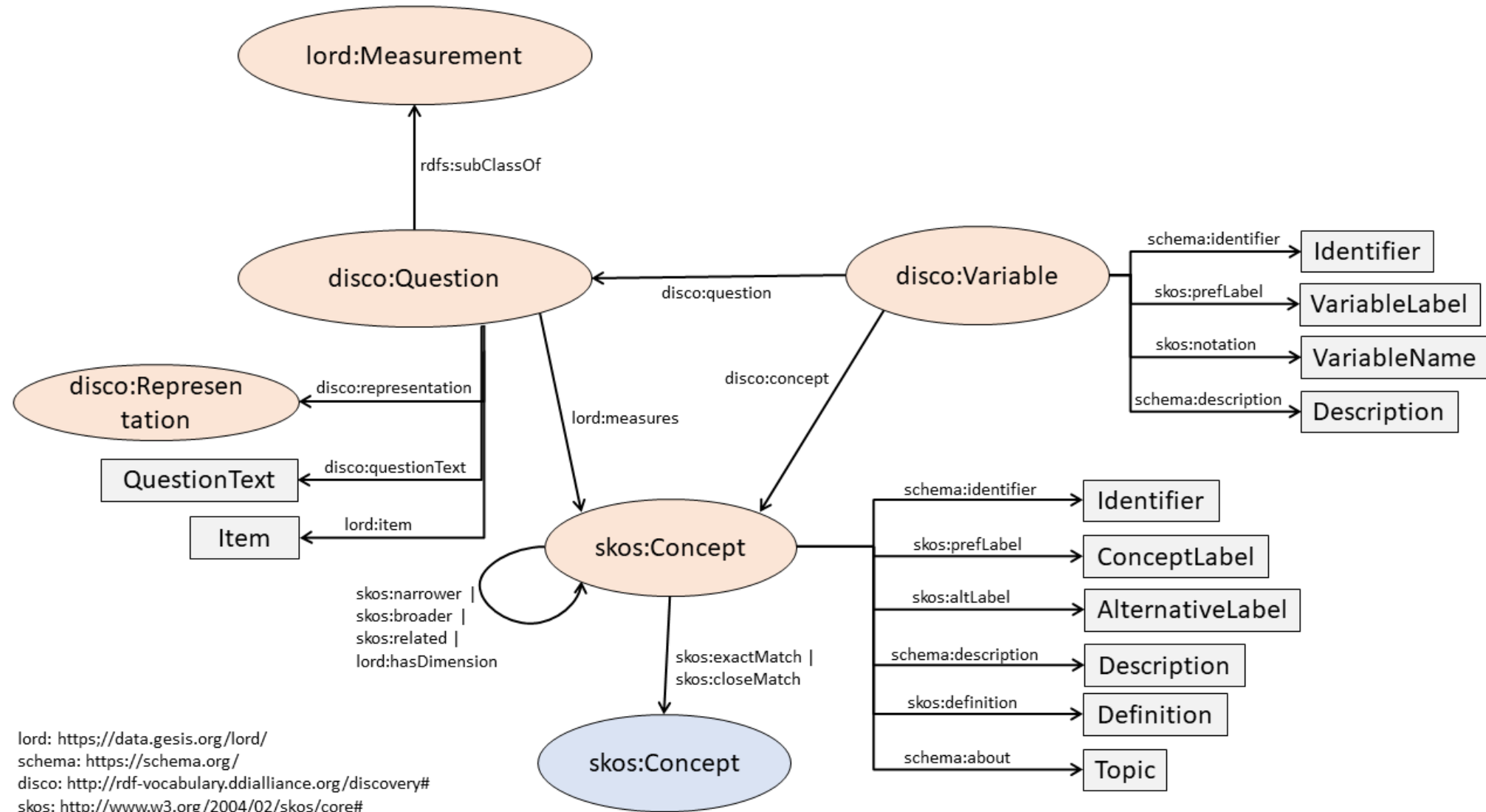
Thank you for your attention

https://www.diw.de/de/diw_01.c.862891.de/projekte/linked_open_research_data_for_social_science_pilot_study_lord_pilot.html

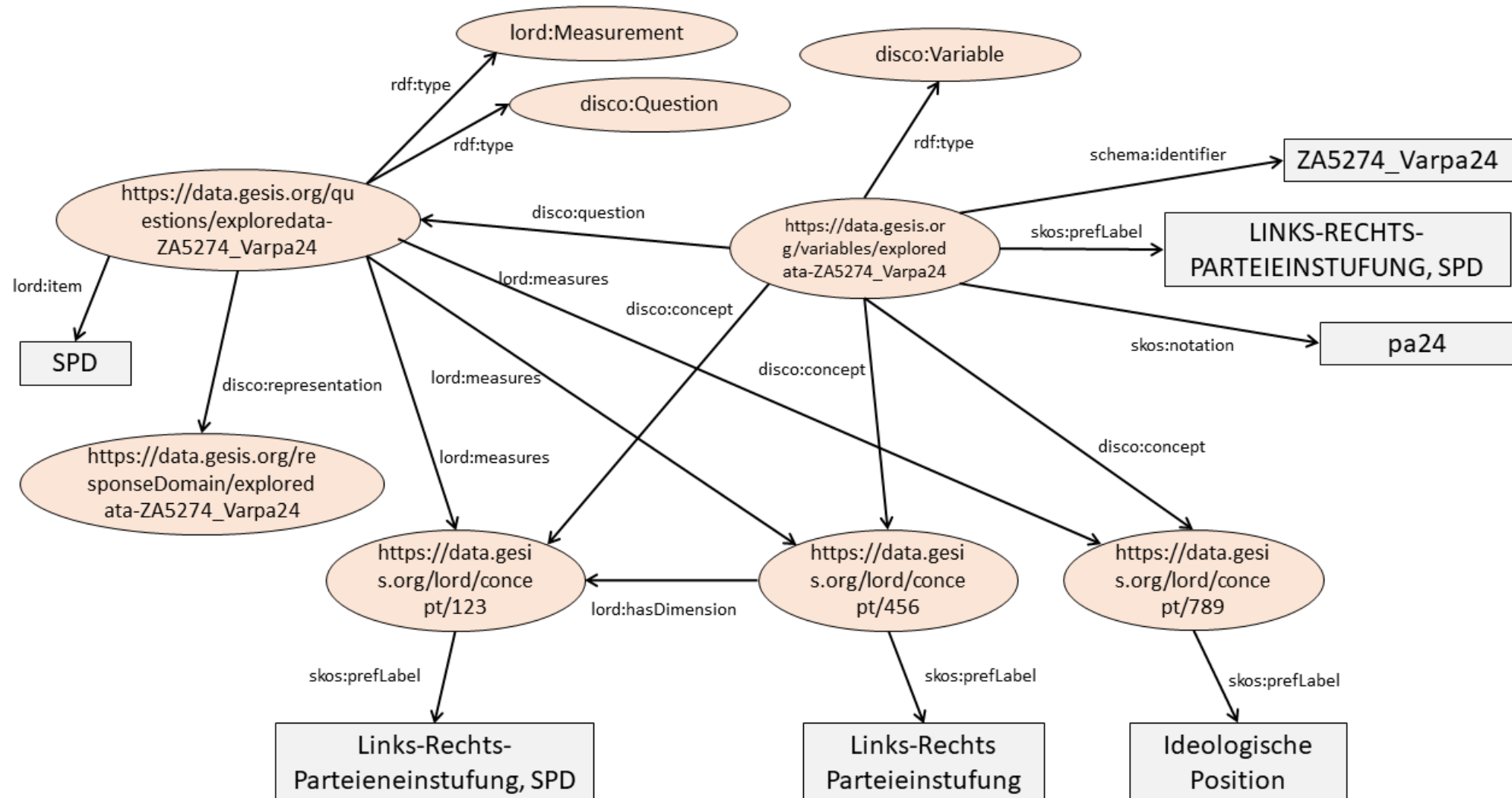


LORDpilot received funding from the German Science Foundation (Grant Number: 464413245)

III. The LORD Data Model (back up)



III. The LORD Data Model: Example (in German)



IV. Lessons Learned from the pilot study (back-up)

