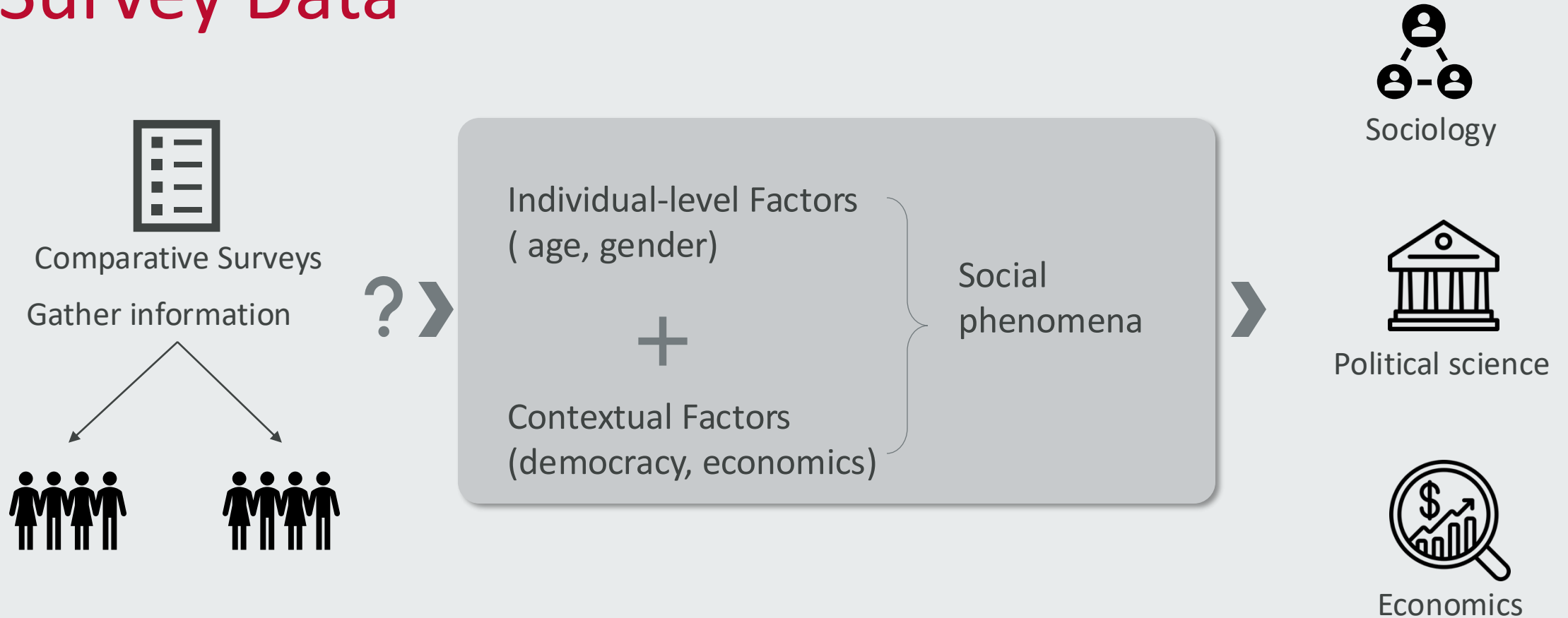# *SDR*Querier: A Visual Querying Framework for Cross-National Survey Data Recycling

Yamei Tu, Olga Li, Junpeng Wang, Han-Wei Shen, Przemek Powałko, Irina Tomescu-Dubrow, Kazimierz M. Slomczynski, Spyros Blanas,J. Craig Jenkins
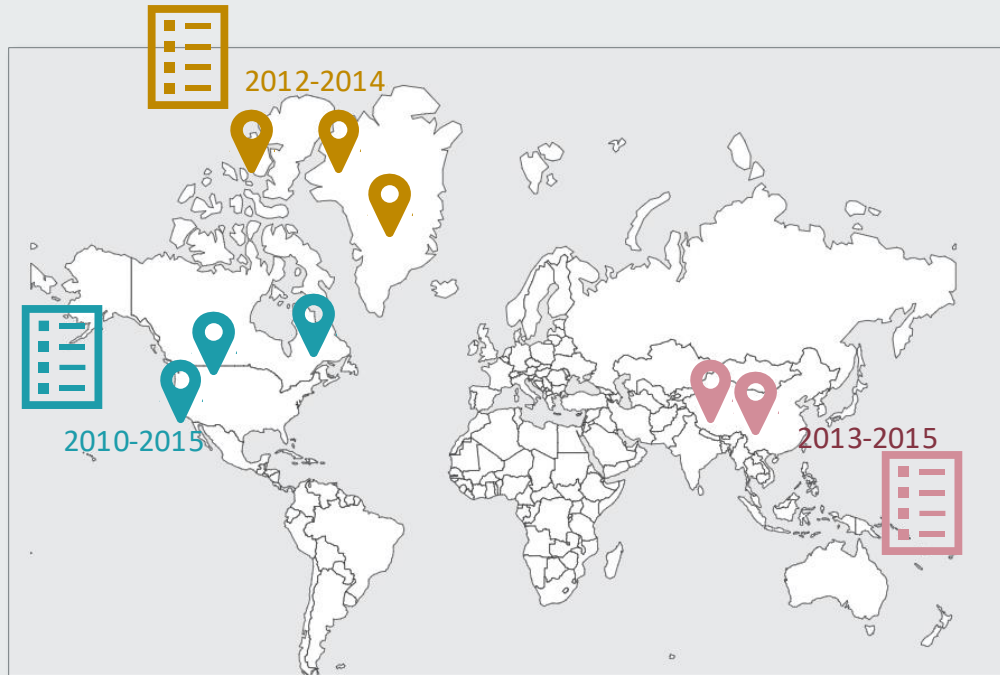
THE OHIO STATE UNIVERSITY

Institute of Philosophy and Sociology, Polish Academy of Sciences (IFiS PAN)

1

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

INSTYTUT FILOZOFII I SOCJOLOGII
POLSKIEJ AKADEMII NAUK

# Survey Data



Comparative Surveys

Gather information

Individual-level Factors ( age, gender)

+

Contextual Factors (democracy, economics)

Social phenomena

Sociology

Political science

Economics

# Survey Data Recycling



2012-2014

2010-2015

2013-2015

22 International Surveys

rescale

recode

transform

Harmonized Dataset

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

INSTYTUT
FILOZOFII
I SOCJOLOGII
POLSKIEJ AKADEMII NAUK
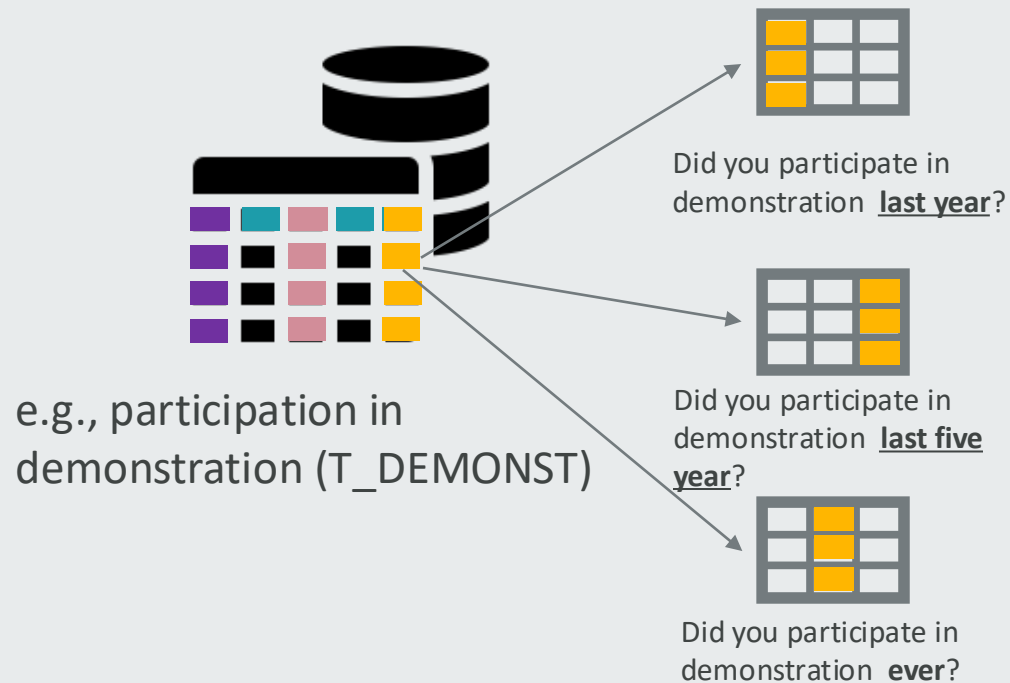
# Harmonized Data

Each row indicates one respondent's responses to all questions.

Each column is called a **variable.**



e.g., participation in demonstration (T_DEMONST)

Did you participate in demonstration **last year**?

Did you participate in demonstration **last five year**?

Did you participate in demonstration **ever**?

**Target Variable**: represents a specific question from the original questionnaire.
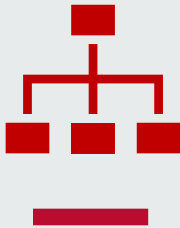
**Source Variable**: original variables taken from different surveys for harmonization

**Harmonization Controls:** inter-survey methodological variability in the formulation of the source questions.

**Quality Controls:** biases and errors that stem from differences in the quality of the source survey data.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

IFiS
INSTYTUT
FILOZOFII
I SOCJOLOGII
P A N
POLSKIEJ AKADEMII NAUK

# Utilizing Harmonized Dataset

Study "gender differences in political participation"?

## Data Understanding

To decide whether there is relevant variables to scientists' research interests?

## Data Exploration

To decide whether there is sufficient data related to target variables.

## Data Evaluation

Run the statistical models on the valid data to verify and generate insights.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

IF IS
INSTYTUT
FILOZOFII
I SOCJOLOGII
P A N
POLSKIEJ AKADEMII NAUK

# Visual Querying Framework: SDRQuerier



## Query-by-Question
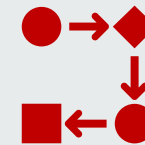
**Whether there are related variables?**

- BERT-based model for variable recommendation from user question.
- Visualizing harmonized dataset structure.



## Query-by-Condition

**Whether variables have sufficient data to use?**

Temporal Availability Profiler to display multi-faceted information



## Query-by-Relation

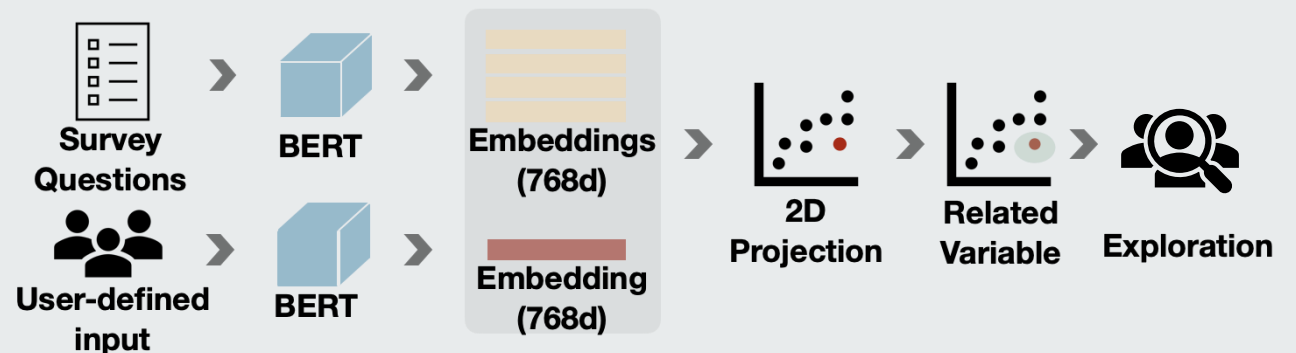**Whether expected patterns exist in the data?**

Presenting relational patterns.

# Visual Querying Framework: SDRQuerier

### Query-by-Question

**Whether there are related variables?**

- BERT-based model for variable recommendation from user question.
- Visualizing harmonized dataset structure.

### Query-by-Condition

**Whether variables have sufficient data to use?**

Temporal Availability Profiler to display multi-faceted information

### Query-by-Relation

**Whether expected patterns exist in the data?**

Presenting relational patterns.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

INSTYTUT
FILOZOFII
I SOCJOLOGII
POLSKIEJ AKADEMII NAUK

# Difficult to identify relevant variables?

- It can be difficult to identify the theoretical concept and their concrete meanings from the abbreviated names.
- It is hard for scientists to find related control variables efficiently due to the complex structure.
- Survey questions can provide good contextual meanings for each target variable.
- e.g., T_DEMONST (a target variable) can be described differently in the original questionnaires:
    - *authorized demonstrations in democratic countries*
    - *unauthorized activities in non-democratic countries*

8

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

IFiS
INSTYTUT
FILOZOFII
I SOCJOLOGII
P A N    POLSKIEJ AKADEMII NAUK

# BERT-based Model

- Developed BERT-based classification model trained on survey questions for predicting target variable.

- Model enables two types of variable recommendation:

  - Hard: predict relevant target variables directly from user input.

  - Soft: generate embeddings of user input and compare to survey question embeddings for more flexible explorations.
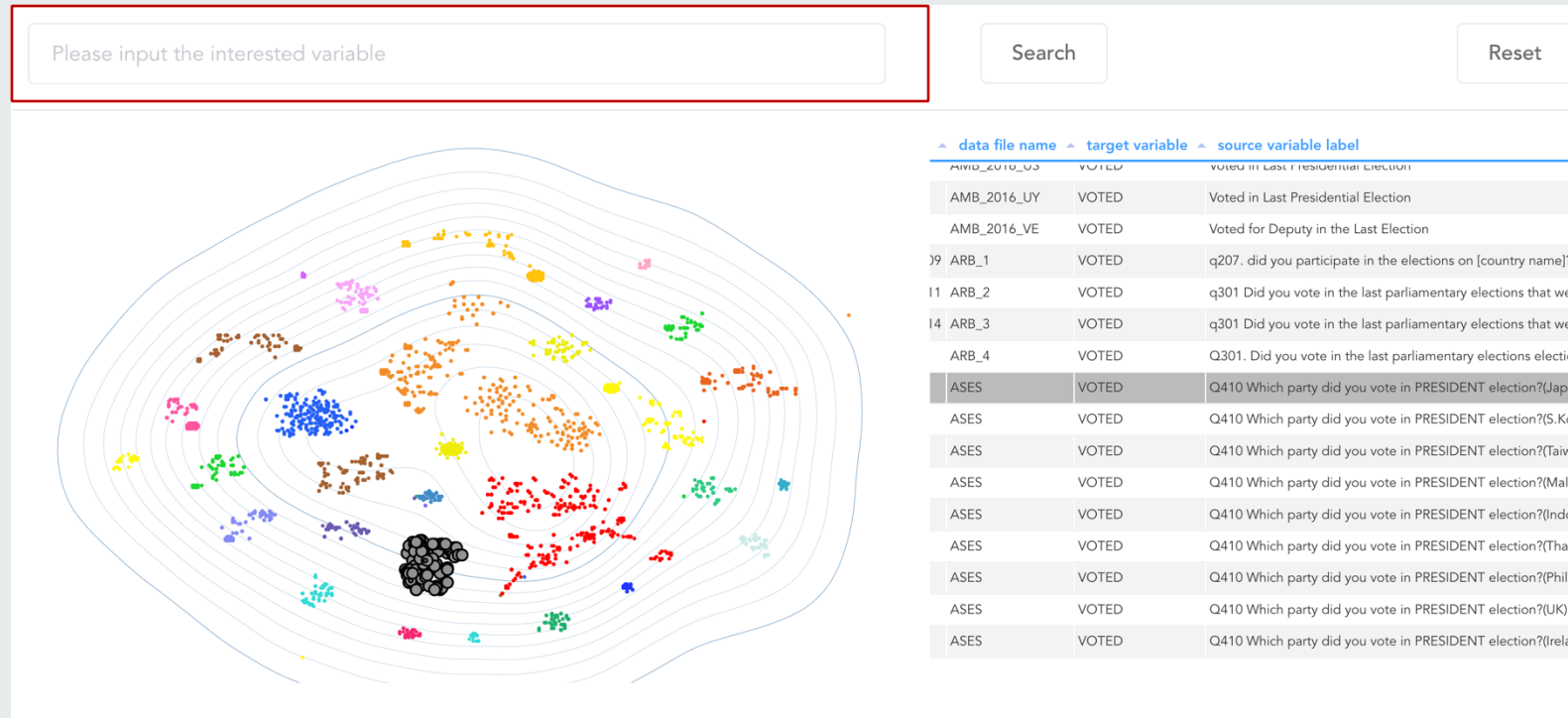
THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

INSTYTUT
FILOZOFII
I SOCJOLOGII
POLSKIEJ AKADEMII NAUK

# Soft recommendation

NLP Model + Visualization + Interactions

# Soft recommendation

## NLP Model + Visualization + Interactions



- Input the questions /keyword /key-phrases to describe the research interests.

THE OHIO STATE UNIVERSITY
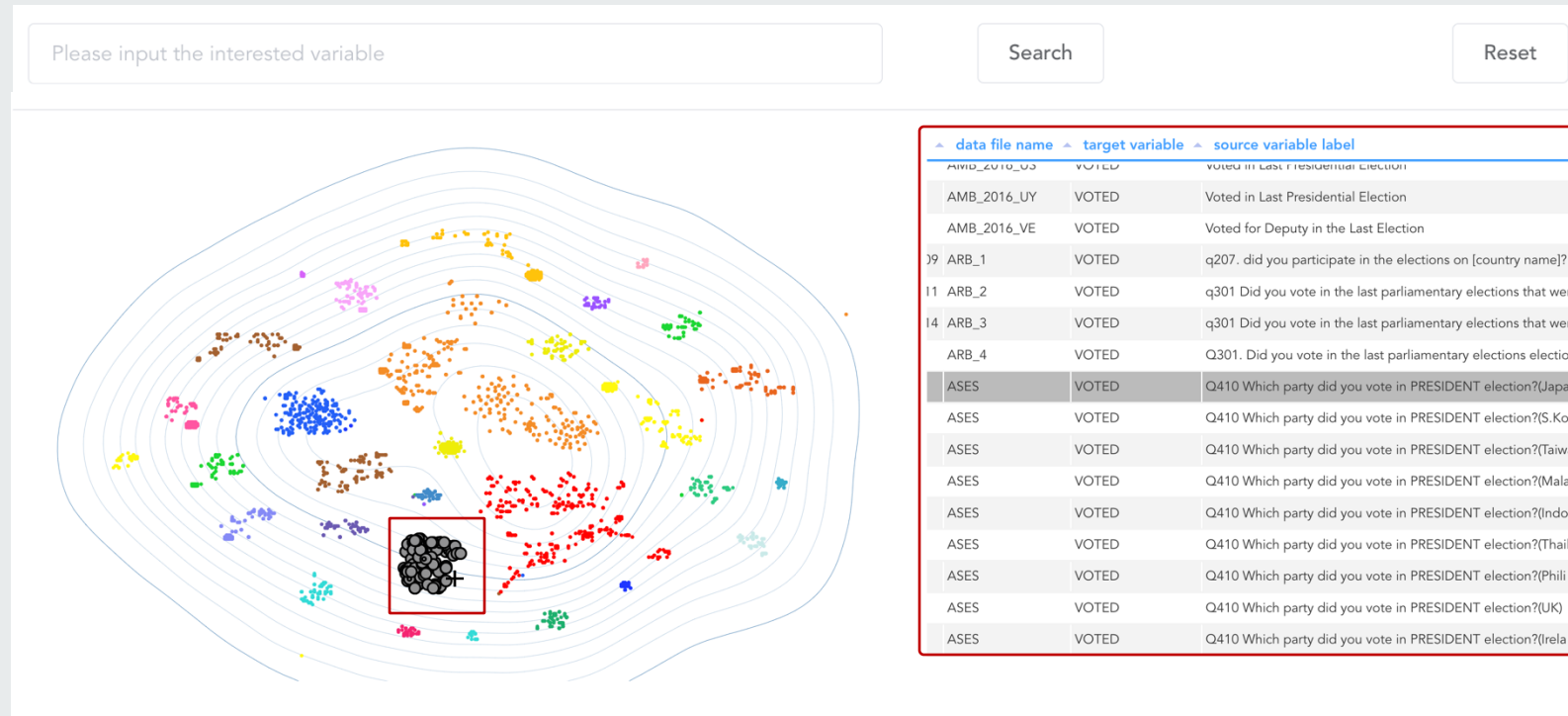COLLEGE OF ENGINEERING
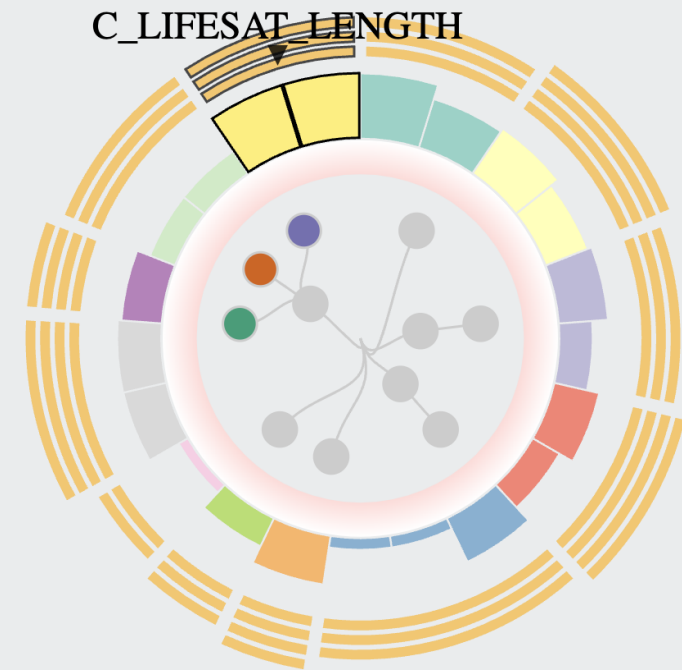
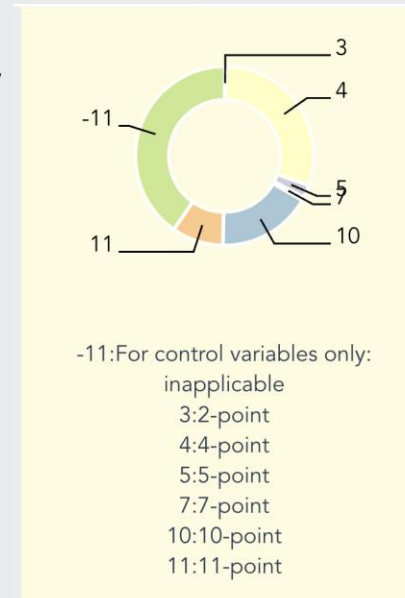# Soft recommendation

## NLP Model + Visualization + Interactions



- Input the questions /keyword /key-phrases to describe the research interests.

- Input embedding from BERT model is projected into the embedding space.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

# Soft recommendation

## NLP Model + Visualization + Interactions



- Input the questions /keyword /key-phrases to describe the research interests.

- Input embedding from BERT model is projected into the embedding space.

- Brushing related questions to understand variable's contextual meanings.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

# Visualizing Harmonized Data Structure

- Harmonized data's structure can be challenging for users to quickly comprehend.
- We design a Circular Graph to
  - Indicate different types of variables
  - Illustrate the relationships among variables.
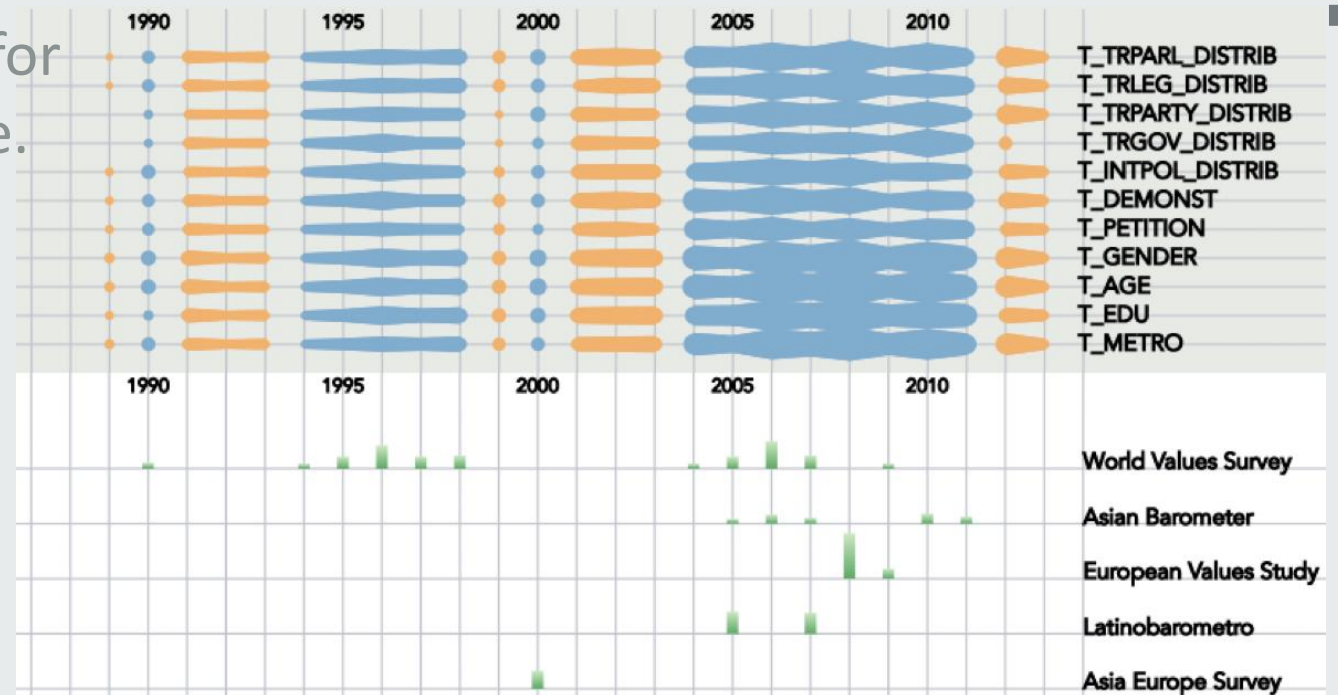- Hard recommendation will highlight the corresponding bars.



-11:For control variables only:
inapplicable
3:2-point
4:4-point
5:5-point
7:7-point
10:10-point
11:11-point

C_LIFESAT_LENGTH

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

INSTYTUT
FILOZOFII
I SOCJOLOGII
POLSKIEJ AKADEMII NAUK

# *Inefficiency in Existing Data Portal*

- Social scientists conduct their study based on a **set of target variables**
- Existing data portal:
  - Blindly download the full harmonized data
  - Ex-post check whether the set of target variables has sufficient data in terms of country and time coverage.
- We design a *Temporal Availability Profiler to* **effectively** and **efficiently** check data availability from multiple perspectives before downloading.
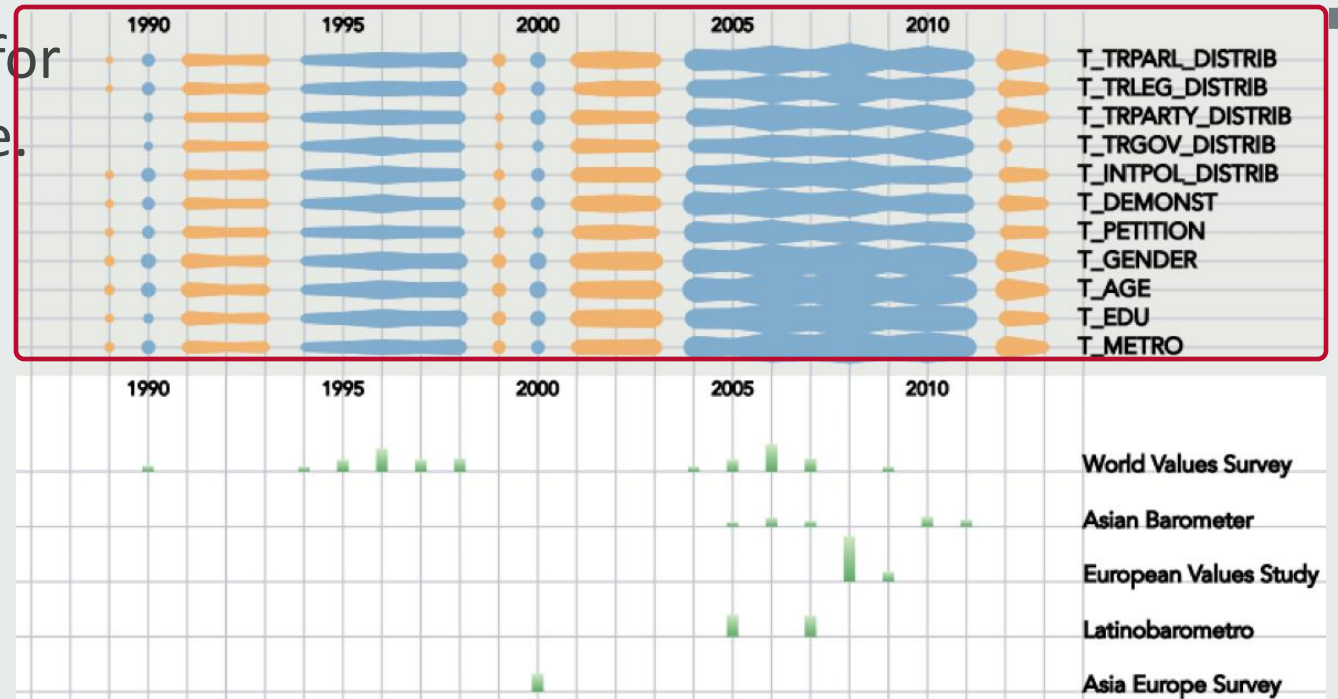
THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

IF iS INSTYTUT FILOZOFII I SOCJOLOGII
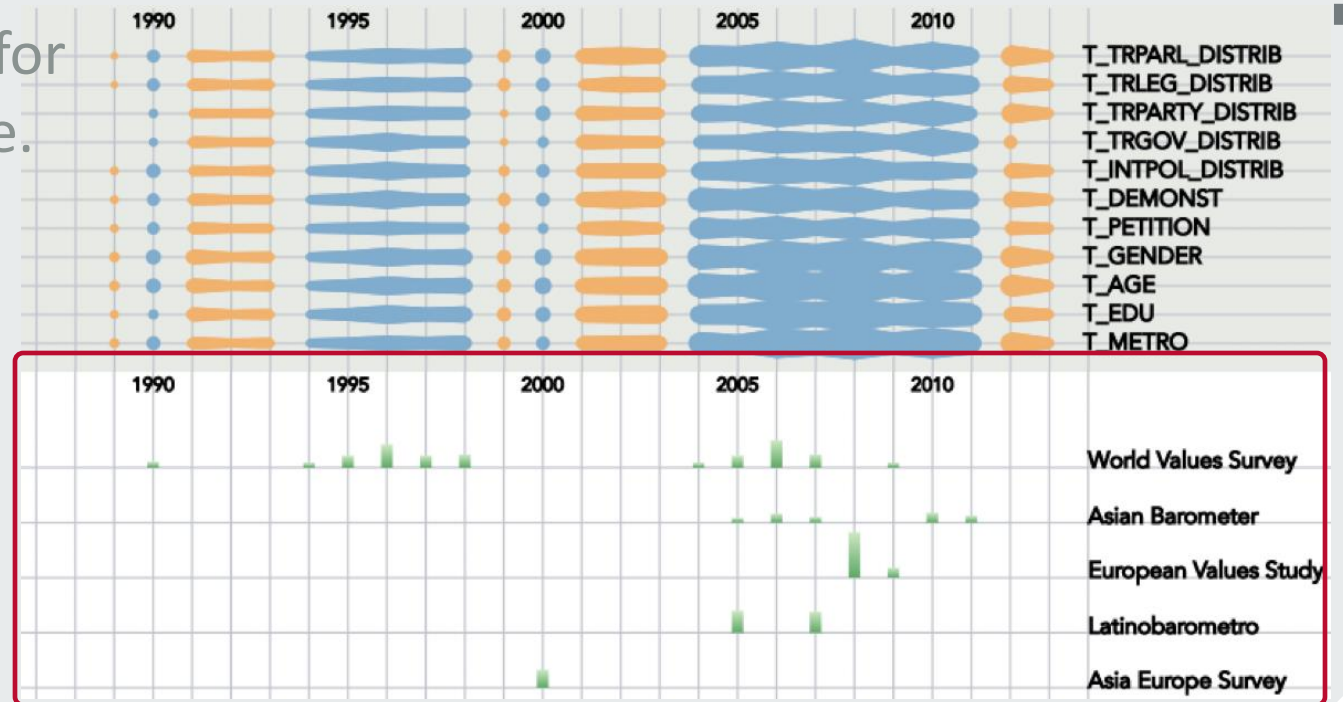P A N POLSKIEJ AKADEMII NAUK

# Temporal Availability Profiler

- It reveals the availability of the harmonized data at multiple levels

- Separate Availability View:
  - the number of valid samples for **each** target variable over time.
  - decide among multiple alternative variables.
- Joint Availability View:
  - available samples for **all** selected target variables.
  - precise pool of valid samples that can be used.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

INSTYTUT FILOZOFII I SOCJOLOGII
POLSKIEJ AKADEMII NAUK
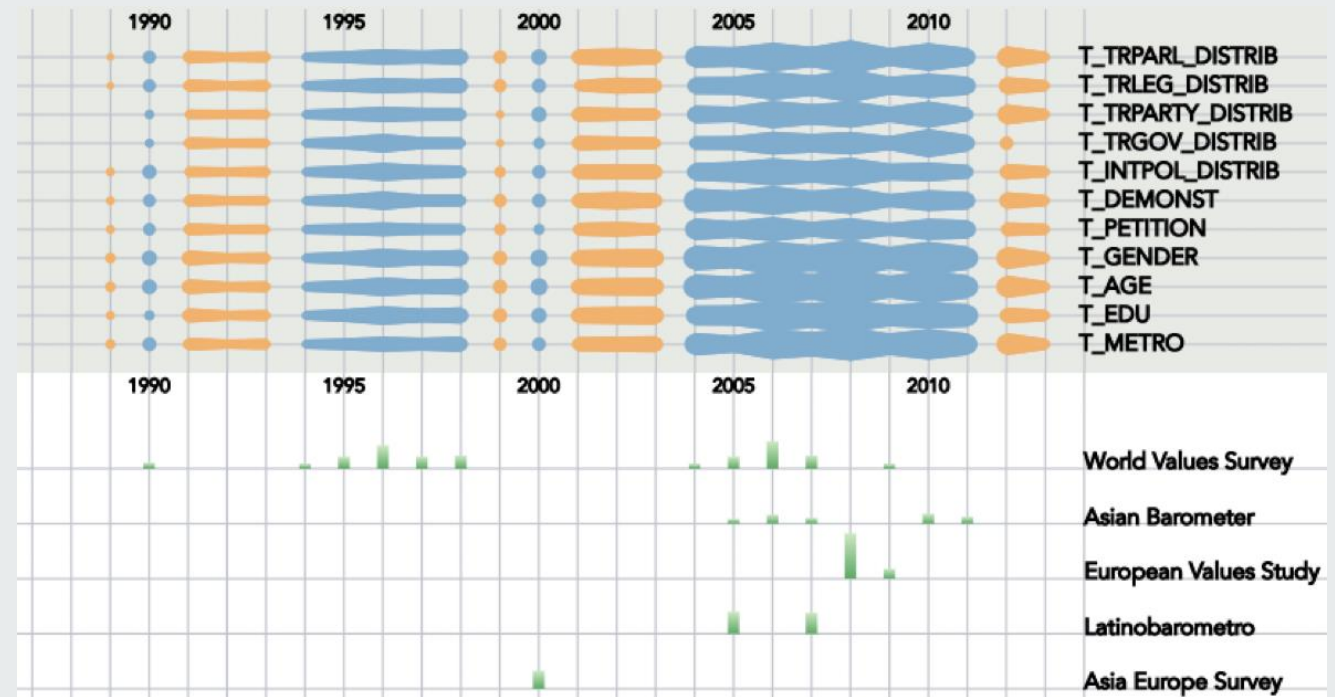
# Temporal Availability Profiler

- It reveals the availability of the harmonized data at multiple levels

- Separate Availability View:
  - the number of valid samples for **each** target variable over time.
  - decide among multiple alternative variables.
- Joint Availability View:
  - available samples for **all** selected target variables.
  - precise pool of valid samples that can be used.

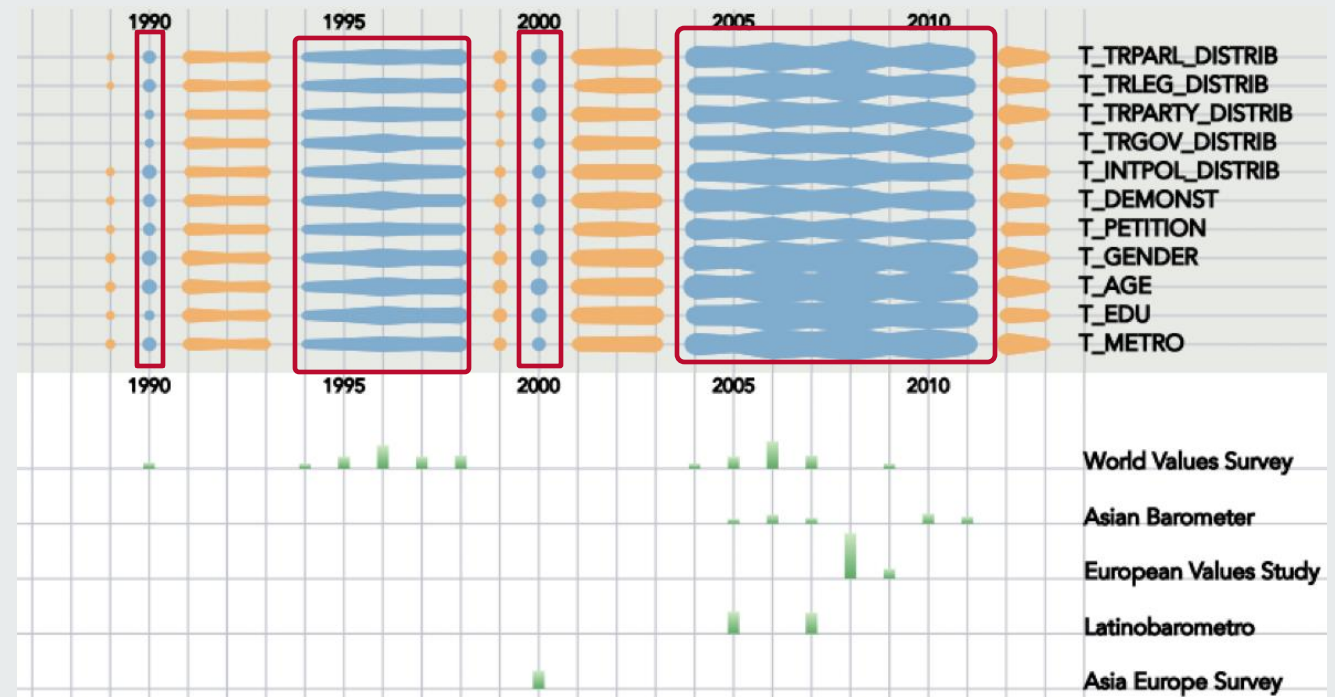# Temporal Availability Profiler

- It reveals the availability of the harmonized data at multiple levels
- Separate Availability View:
  - the number of valid samples for **each** target variable over time.
  - decide among multiple alternative variables.
- Joint Availability View:
  - available samples for **all** selected target variables.
  - precise pool of valid samples that can be used.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

INSTYTUT FILOZOFII I SOCJOLOGII
POLSKIEJ AKADEMII NAUK
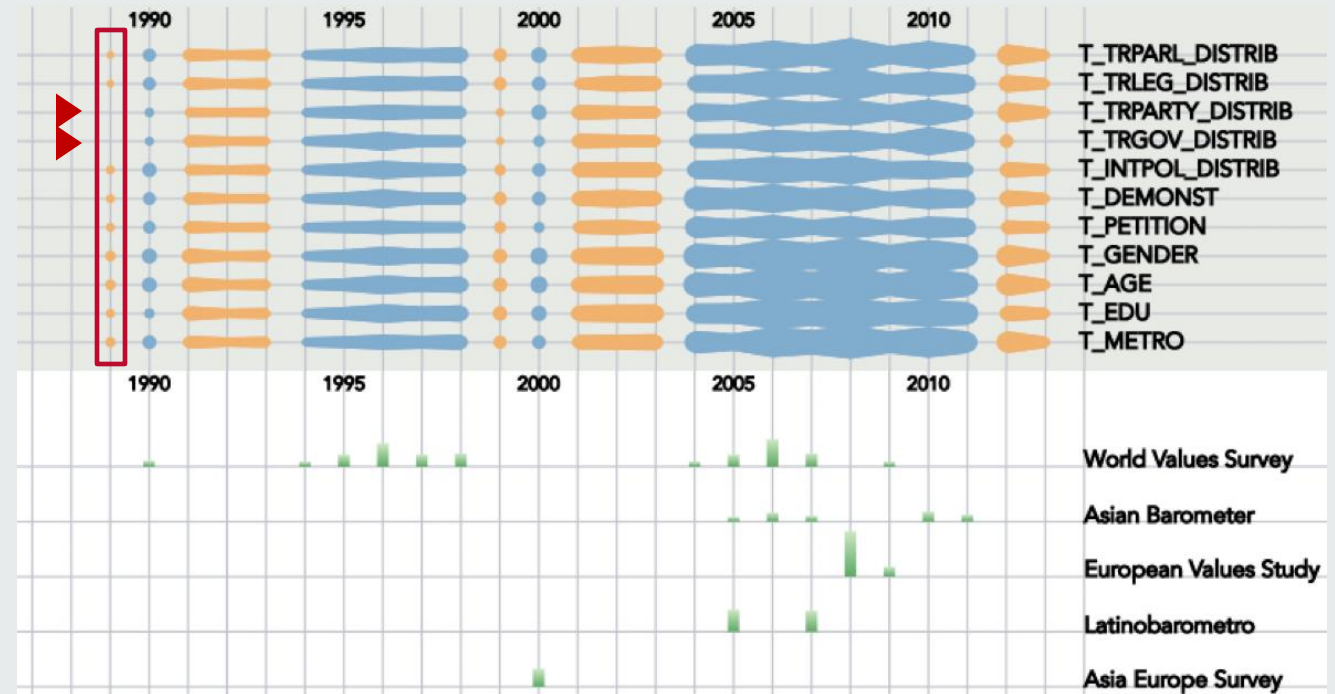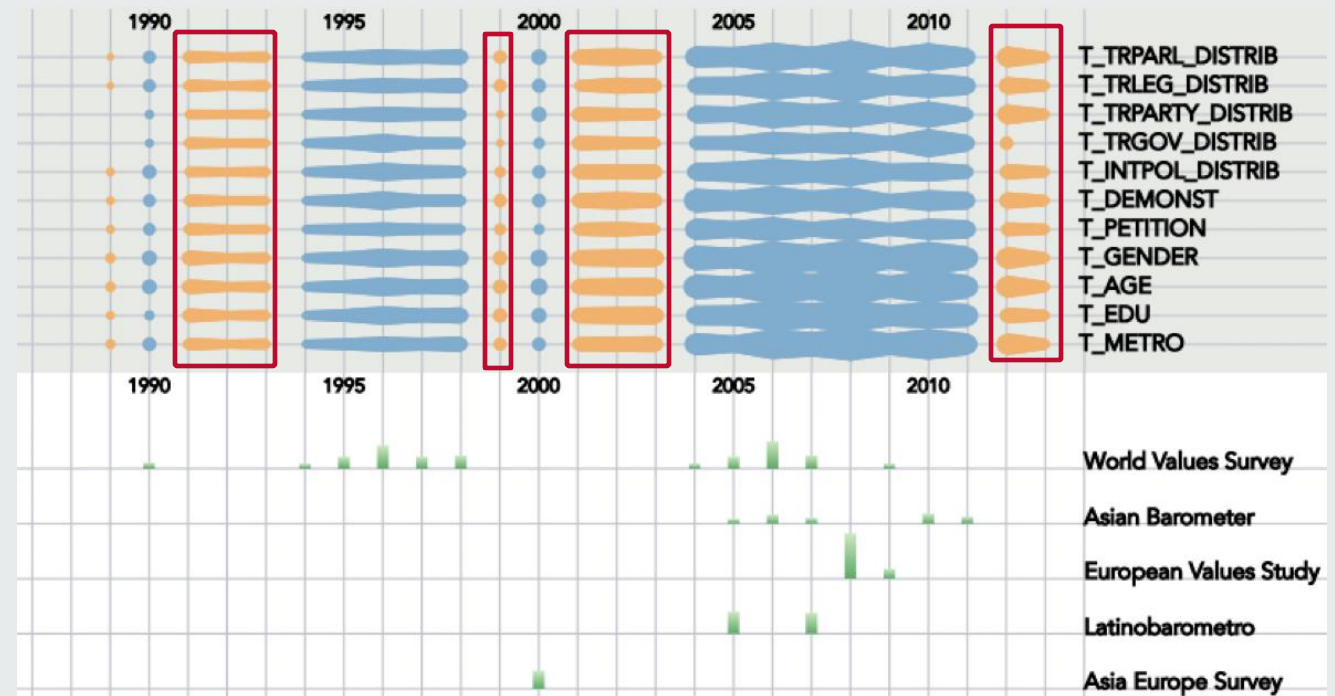
# Temporal Availability Profiler

Three cases:

1. Each target variable has data, and there are also jointly available samples.

2. At least one target variable does not have data, resulting in the lack of jointly available data.

3. Each target variable has data, but there is no overlap among them.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

# Temporal Availability Profiler

Three cases:

1. Each target variable has data, and there are also jointly available samples.

2. At least one target variable does not have data, resulting in the lack of jointly available data.

3. Each target variable has data, but there is no overlap among them.

# Temporal Availability Profiler

Three cases:

1. Each target variable has data, and there are also jointly available samples.

2. At least one target variable does not have data, resulting in the lack of jointly available data.

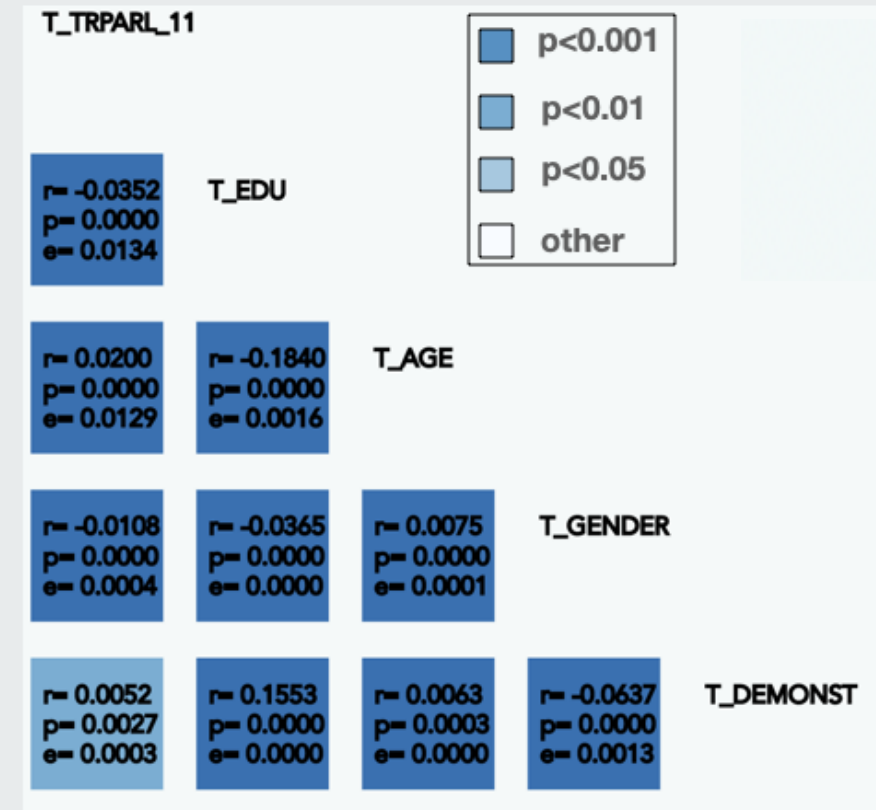3. Each target variable has data, but there is no overlap among them.

# Temporal Availability Profiler

Three cases:

1. Each target variable has data, and there are also jointly available samples.

2. At least one target variable does not have data, resulting in the lack of jointly available data.

3. Each target variable has data, but there is no overlap among them.

# Query-by-Relation

- Social Scientists proposes ***theoretically derived hypotheses***, which are tested through ***statistical models*** using appropriate data.

- QBR is designed to query the hidden patterns from data for model verification and improvements.

- It answers two questions:

  - *what are the relationships between the selected target variables?*

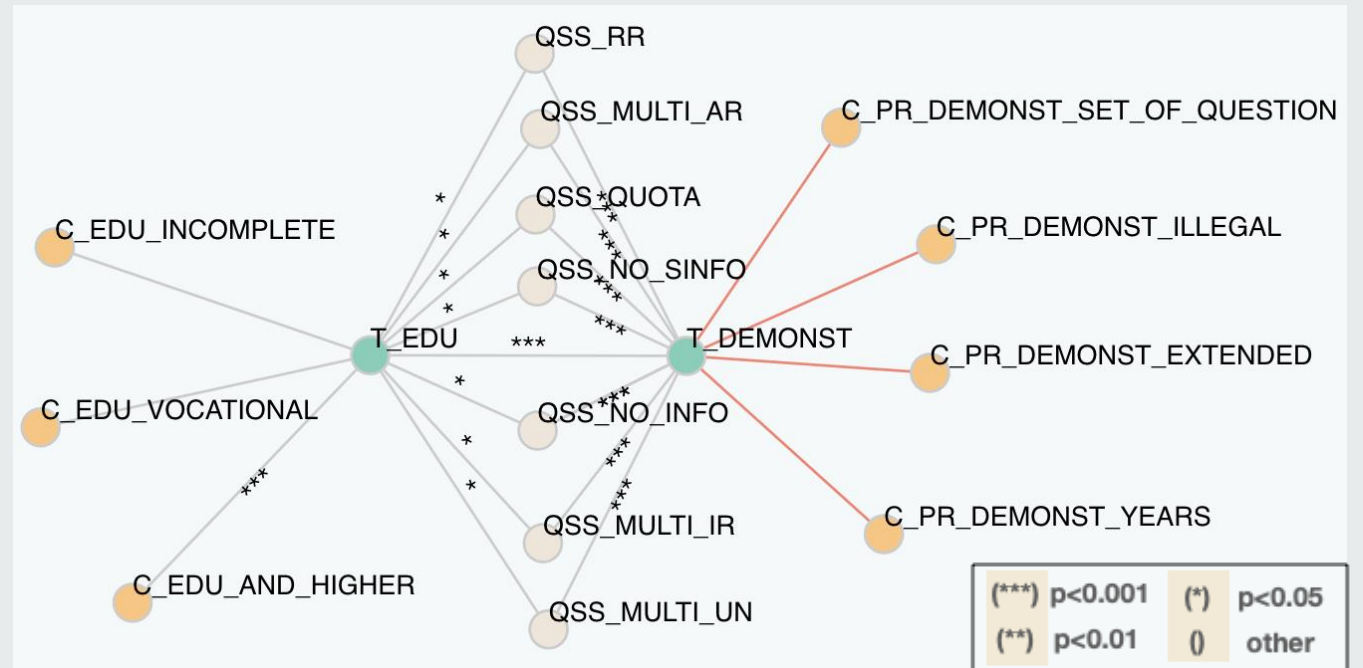  - *what are the potentially related variables?*

# Correlation Matrix

- Showing pairwise correlations for user-selected target variables.
- Helps users to determine what to keep for their regression analysis.
- We compute:
  - Pearson correlation coefficient
  - p-values,
  - standard errors.

THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

# Network Visualization

- Query complex relations given one pair of target variables.
- Both types of methodological variables can affect the relationship between substantive variables
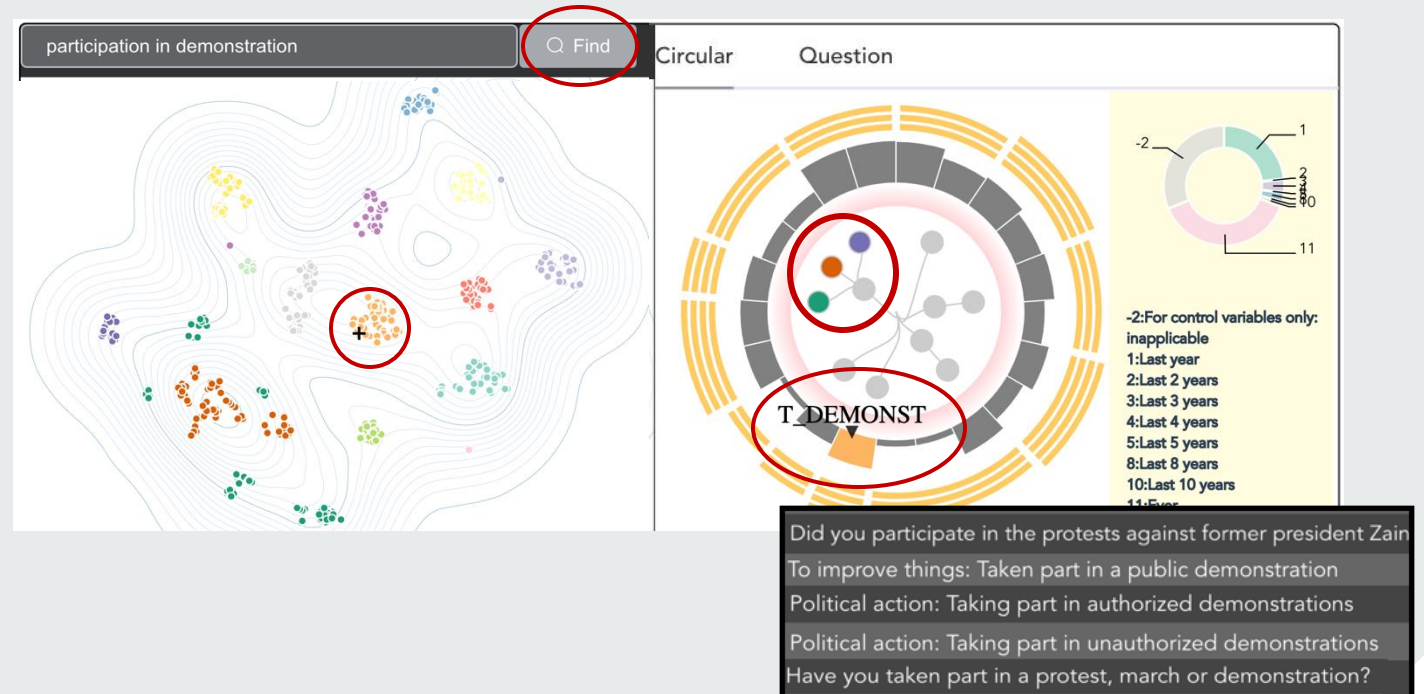  - Control variable
  - Quality variable

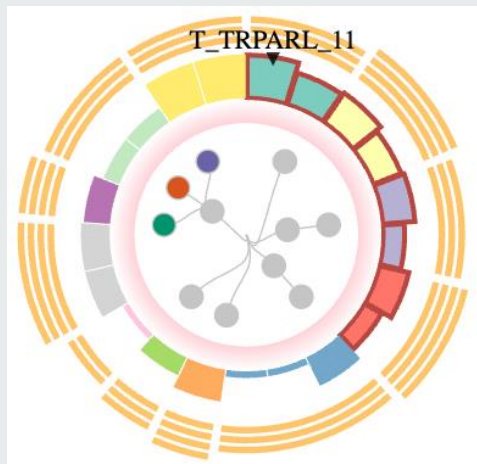THE OHIO STATE UNIVERSITY
COLLEGE OF ENGINEERING

# Case Study: Participation in Demonstrations worldwide

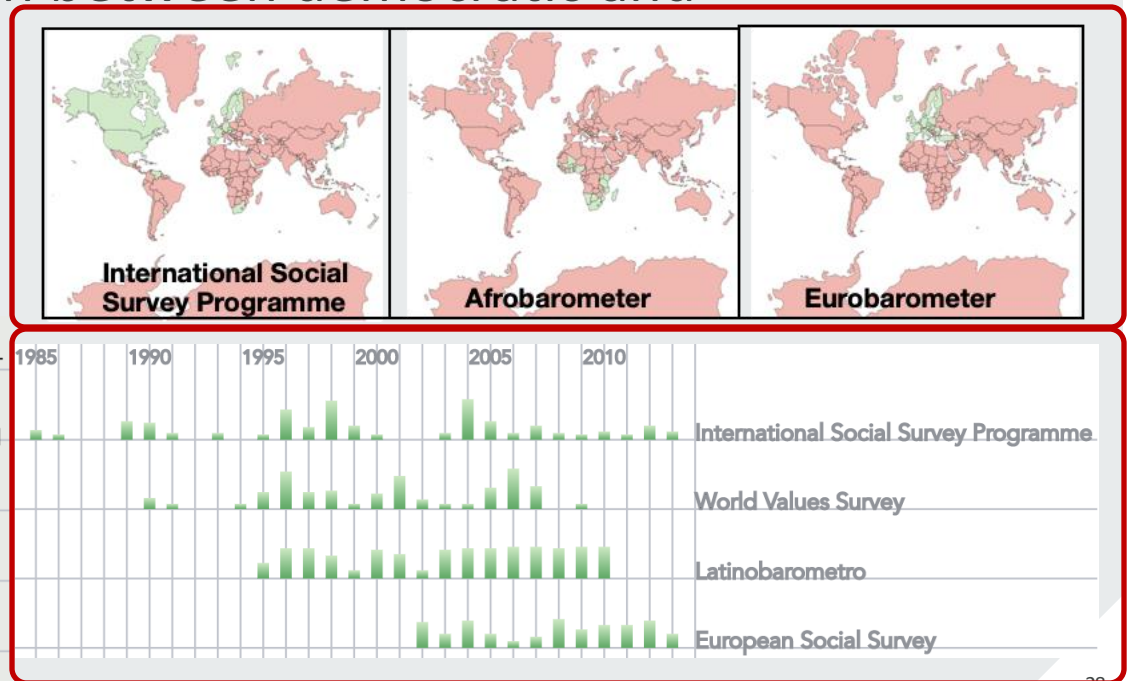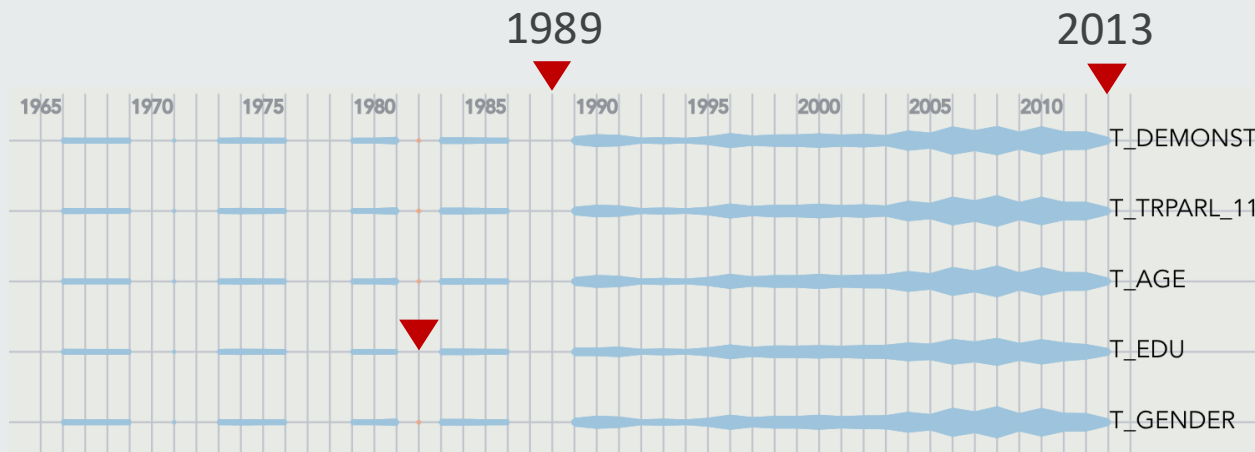T_EDU, T_AGE, T_GENDER          T_TRPARL_11

- Proposal: that resources and political attitudes have different effects on the levels of participation in demonstrations in democratic and non-democratic countries.

T_DEMONST

# Case Study: Participation in Demonstrations worldwide

- High Temporal Coverage
- Harmonized data provides 22 available surveys to use
- Diverse country coverage allows comparison between democratic and non-democratic countries.

# Conclusion

- We abstract the challenges in analyzing harmonization survey data
- We propose **a visual analytics system**, *SDR*Querier, to help scientists locate target variables and evaluate their theoretical models.

**THE OHIO STATE UNIVERSITY**

COLLEGE OF ENGINEERING

# Acknowledgement

# Questions?

THE OHIO STATE UNIVERSITY

COLLEGE OF ENGINEERING