



CROSS-NATIONAL  
DATA CENTER  
*in Luxembourg*

# Income Database of Harmonized Microdata: LIS- Cross National Data Center in Luxembourg

Think thanks Summit  
Zurich, 22 January 2019

# Outline of the presentation

- I. LIS, the institution
- II. The LIS/LWS Databases
- III. Research possibilities
- IV. Data dissemination & Documentation
- V. LIS golden rules for harmonization



# Part I:

# *LIS, the institution*



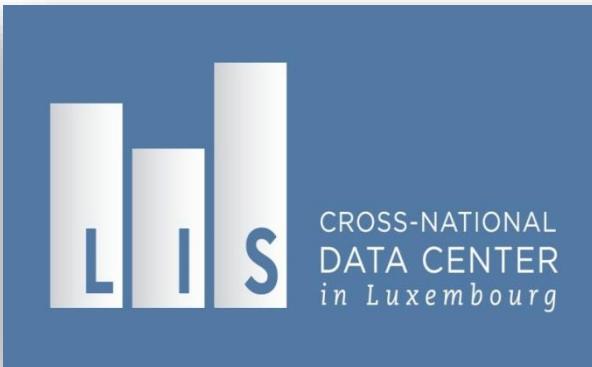
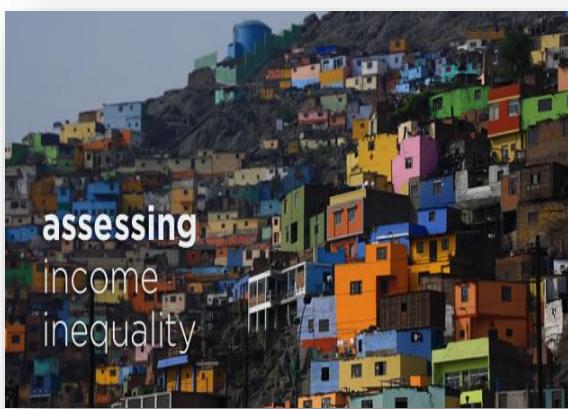
# LIS History

- LIS was founded in 1983 by two US academics (Professors Tim Smeeding and Lee Rainwater) and a team of multi-disciplinary researchers in Europe. It began as a “study”, which later grew and was institutionalized as “LIS”. (Note the “oral history” video on LIS website.)
- For nearly 20 years, LIS was part of a Luxembourg-based research institute known as CEPS (recently renamed LISER). In 2002, LIS became an independent non-profit institution (an ASBL).
- LIS is supported by the Luxembourg government (25%), by the national science foundations and other funders in many of the participating countries (50%), and by several supranational organizations – supplemented by project-specific grants and some private philanthropy (25%).



# LIS Mission

To enable, facilitate, promote, and conduct cross-national comparative research on socio-economic outcomes and on the institutional factors that shape those outcomes.



# LIS: new leadership structure

launched 1 September, 2016

**Prof. Daniele Checchi**

**Director of Luxembourg Office of LIS**

\* He is Professor of Economics currently on leave from the University of Milan.

\* He is currently working at the Italian National agency for the Evaluation of the University system.



**Prof. Janet Gornick**

**Director of US Office of LIS**

\* She is Professor of Political Science and Sociology, at the Graduate Center of the City University of New York.

\* The “LIS Center” has been renamed “The Stone Center on Socio-Economic Inequality”.

\* The “US Office of LIS” is now an entity within the enlarged Center.



# LIS structure

## LIS - Cross-National Data Center

François Bourguignon, President

LIS Executive Committee

Daniele Checchi, Secretary General

### MAIN OFFICE:

#### Luxembourg Office of LIS

Daniele Checchi, Director

### SATELLITE OFFICE:

#### US Office of LIS

Janet Gornick, Director

- parent organization
- located in Luxembourg, together with the University and LISER
- independent, chartered non-profit organization
- cross-national, participatory governance
- acquires, harmonizes, and disseminates data for research
- venue for research, visiting scholars, conferences, and user training

- satellite office
- located at the Graduate Center of the City University of New York
- one pillar of *Stone Center on Socio-Economic Inequality*
- administrative, managerial, development support to parent office
- collaborative public programs (lectures, conversations)
- venue for research, teaching, and graduate student supervision

### LIS Senior Scholars



# LIS: who's who?

## LIS - Cross-National Data Center

François Bourguignon, President

LIS Executive Committee

Daniele Checchi, Secretary General

MAIN OFFICE: Luxembourg Office of LIS Daniele Checchi, Director	SATELLITE OFFICE: US Office of LIS Janet Gornick, Director
<p>Thierry Kruten, Director of Operations and IT Lucie Scapoli, Administrator Officer Benjamin Gérard, System and Network Administrator Data team: Teresa Munzi, Data Team Manager Paul Alkemade              Heba Omar Andrey Cupak              Piotr Paradowski Jörg Neugschwender      Carmen Petrovici</p>	<p>Caroline Batzdorf, Assistant Director Mei-Ling Israel, Financial Manager Laurie Maldonado, Senior Administrator</p>

## LIS Senior Scholars

Prof. Louis Chauvel  
Prof. Daniele Checchi  
Prof. Conchita D'Ambrosio  
Prof. Markus Jäntti  
Prof. Frank Cowell

Prof. Janet Gornick  
Prof. Paul Krugman  
Prof. Leslie McCall  
Prof. Branko Milanovic



# LIS' partners

Our partners include data providers, data users, and funders, in more than 50 countries ...  
and in major supranational organizations, including:

## Financial contributors:

World Bank (WB)

Organization for Economic Cooperation and Development (OECD)

International Monetary Fund (IMF)

International Labor Organization (ILO)

European Union (InGRID)

## Dataset exchange; joint research projects; joint fundraising:

Economic Research Forum (ERF)

European Central Bank (ECB)

French Development Agency (AFD)

LISER

University of Luxembourg



# Part II: The *LIS/LWS Databases*



# Overview of the LIS data

## Deliverables

Two cross-national harmonised databases that allow international comparative research using micro-data:

- LIS (focus on income): 339 datasets – **this presentation will focus on LIS**
- LWS (focus on wealth): 39 datasets

## Scope

Initial focus on high-income countries, successively extended to middle-income countries

## Time span

From the late 1960s to up-to-date

## Geographical coverage

World-wide, but some regions are less covered (Africa, EECCA...)

## Main contents

- household composition and characteristics
- socio-demographic characteristics of household members
- extensive set of labour market data
- detailed breakdown of household and individual income data
- household consumption data
- a detailed set of wealth and behavioural variables (LWS only)



# Luxembourg Income Study Database (LIS)

- First and largest available database of harmonized income data, available at the household and person levels
- In existence since 1983
- Data mostly start in 1980, some go back to the 1960s (recollected every 3-5 years)
- Started with six countries; now 50 countries
- 300+ datasets (repeated cross sections)
- Used to study: poverty; income inequality; labor market outcomes; policy effects

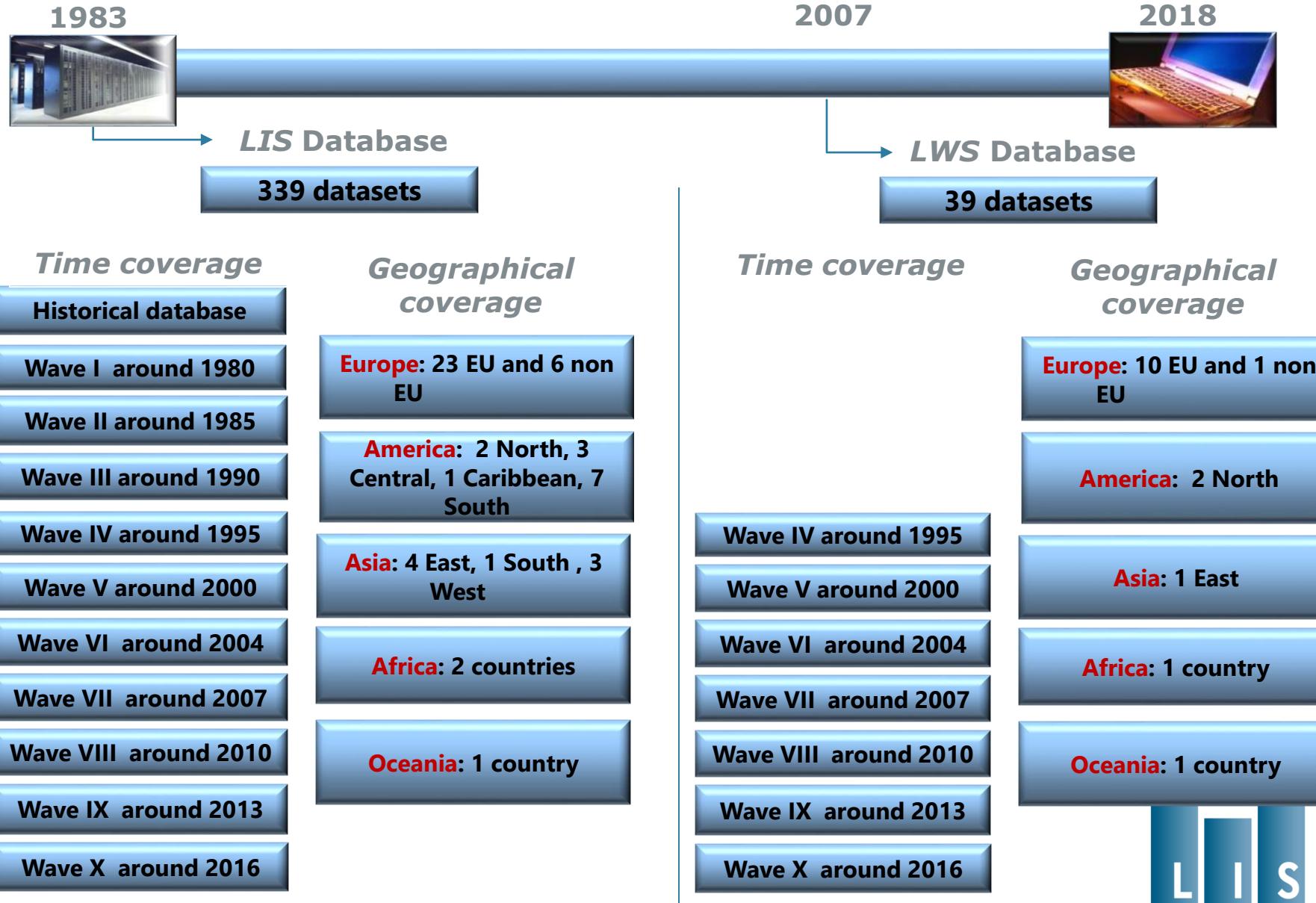


# Luxembourg Wealth Study Database (LWS)

- First available database of harmonized wealth data, available at the household level
- In existence since 2007
- 52 datasets from 16 countries – up or in process
- Revised and updated in 2016 (in coordination with Euro system's Household Finance and Consumption Survey - HFCS)
- Used to study: *household assets, debt, and expenditures; wealth portfolios; policy effects.*



# The LIS Databases: *LIS* and *LWS*



# Countries in LIS and LWS Databases (N=50)

*approximately 65% of world population and 85% of world GDP*

High-income countries (N=34):			Upper-middle-income countries (N=12):		Lower-middle-income countries (N=4):
Australia	Hungary	Slovak Republic	Brazil	Paraguay	Cote D'Ivoire (q1 '19)
Austria	Iceland	Slovenia	China	Peru	Egypt
Belgium	Ireland	South Korea	Colombia	Romania	Georgia
Canada	Israel	Spain	Dom. Republic	Russia	India
Chile	Italy	Sweden	Guatemala	Serbia	
Czech Republic	Japan	Switzerland	Mexico	South Africa	
Denmark	Lithuania	Taiwan			
Estonia	Luxembourg	United Kingdom			
Finland	Netherlands	United States			
France	Norway	Uruguay			
Germany	Panama				
Greece	Poland				

# LIS country coverage by end of 2018



# What we do at LIS

## Core activity: Data work

### Step 1. Data acquisition

We identify appropriate datasets (*reliable, and high-quality data*)

We negotiate with each data provider

### Step 2. Data harmonisation

Common cross-national template

Comprehensive documentation

### Step 3. Data dissemination

We create national-level indicators (*LIS Key Figures*)

We provide harmonized microdata to researchers via *remote execution*

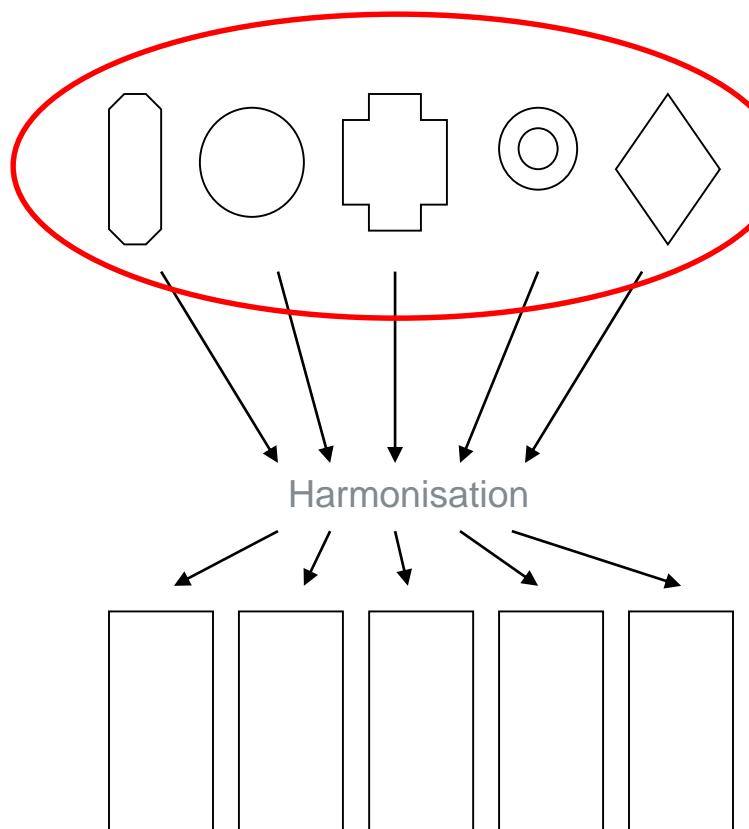
## Other activities

*Research-promotion activities* (conferences, work on methodological issues, collaborations with networks/users/journalists, newsletter, individual research)

Support (user support, research visits)

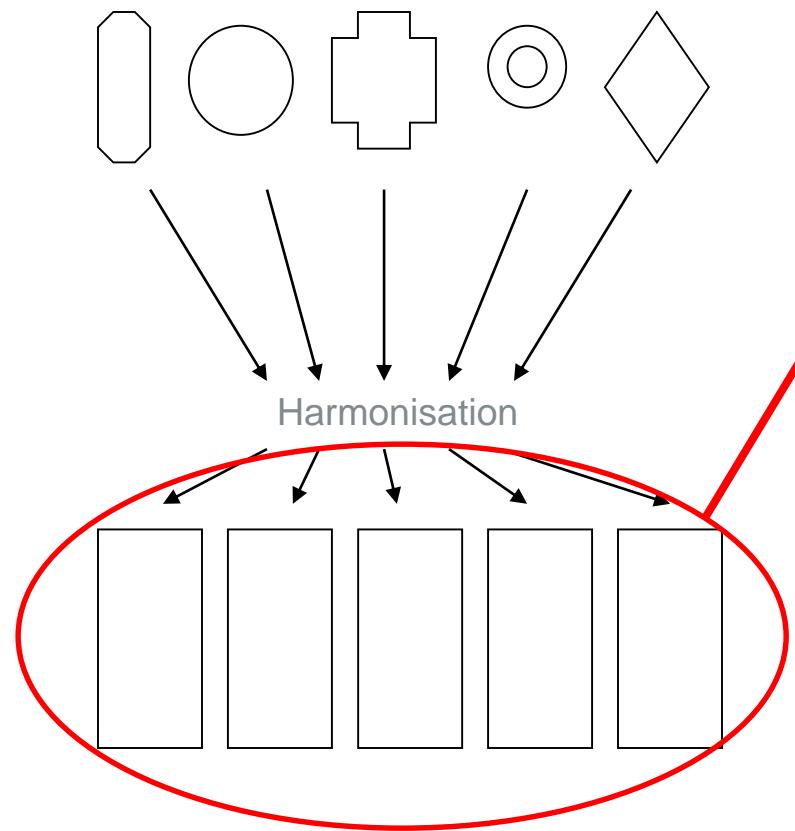


# Ex-post harmonisation at LIS



- The origins of the LIS data
  - ❑ LIS does not organise surveys but collects data from existing data sources:
    - Survey data: *income, household budget, living conditions, multipurpose, human development*
    - Administrative records: *tax records, employers records, social security records*
    - Any mix of the above
  - ❑ Common denominator:
    - *microdata (household and individual level)*
    - *representative of the whole population*
    - *good quality income/wealth data*
    - *main demographic and (possibly) labour market information*

# *Ex-post harmonisation at LIS*



## Final output: the LIS/LWS datasets (CCYY)

- Technical harmonisation: same file structure, same variables

## LIS files

LIS Household File (H)			
HID	...	...	...
LIS Person File (P)			
HID	PID	...	...

## LWS files

LWS Household File (H)

HID	INUM	...	...

LWS Person File (P)

HID	PID	INUM	...

LWS Replicate Weights File (R)

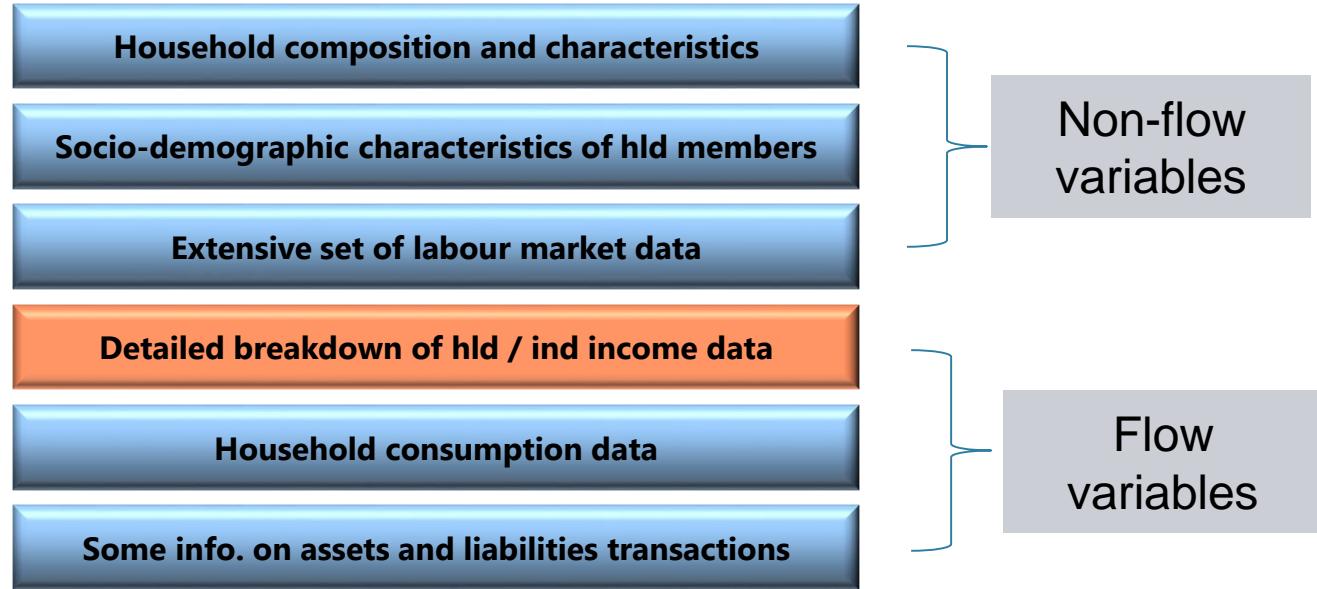
HID	...	...	...

- Conceptual harmonisation
    - Based on the same definitions
    - Comparable concepts

Harmonisation allows LIS users to eliminate many of the potential sources of technical and conceptual non-comparability



# LIS Database content



Household disposable  
income



# LIS Flow variables

## Current income

- Monetary payments as well as (the value of) goods and services received by the household or by individual members of the household at periodic intervals (annual or smaller), that are available for current consumption and that do not reduce the net worth of the household.

## Windfall income

- Windfall gains and other such irregular and typically one-time receipts

## Non-Consumption expenditure

- Monetary expenditures (i.e. paid directly by the household and/or its members) and nonmonetary expenditures (paid on behalf of the household and/or its members) on non-consumption goods and services (such as taxes, contributions, donations, inter-household transfers and interest paid).

## Consumption

- Monetary and non-monetary consumption items.

## Assets/Liabilities transactions

- Monetary inflows that do not constitute income (sales of real estate, financial products, durables or inflows from loans) and outflows that do not represent consumption (purchase of real estate financial products or outflows from loans)

# LIS Household Disposable Income

		MONETARY	NON-MONETARY
LABOUR	Dependent employment	Wages, salaries, bonuses	In-kind earnings
	Self-employment	Profits and losses	Own consumption
+	Financial investment	Interest and dividends	-
CAPITAL	Real estate investment	Rental income	Imputed rent
	Social security transfers		
	- <i>Work-related insurance</i>	Insurance pensions and wage-replacement benefits	-
+ TRANSFERS	- <i>Universal transfers</i>	Universal pensions and universal benefits	STIK
	- <i>Social assistance transfers</i>	Minimum income guarantee Inter-household transfers, transfers from charity	In-kind social assistance In-kind benefits from privates
	Private transfers		
<b>= TOTAL GROSS INCOME</b>			
-		Income taxes	
DEDUCTIONS		Social security contributions	
<b>= HOUSEHOLD DISPOSABLE INCOME</b>			

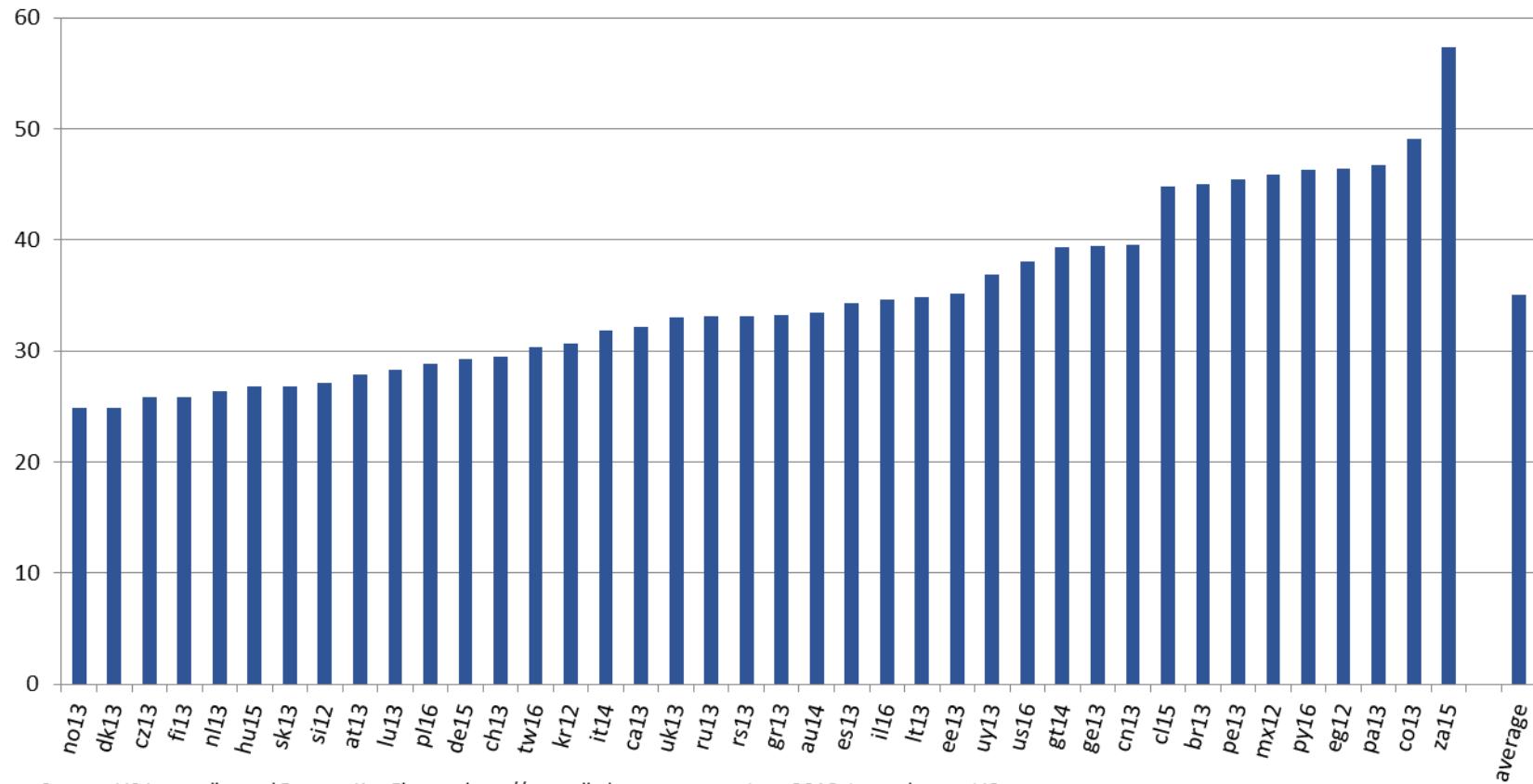


# Part III: Research possibilities



# Measuring Inequality Across Households

Gini Coefficient for LIS countries  
(latest available year)



Source: LIS Inequality and Poverty Key Figures, <http://www.lisdatacenter.org>. June 2018. Luxembourg: LIS.

# Measuring Poverty

## Household Poverty Rates

### Relative poverty rates for the overall population, children and elderly

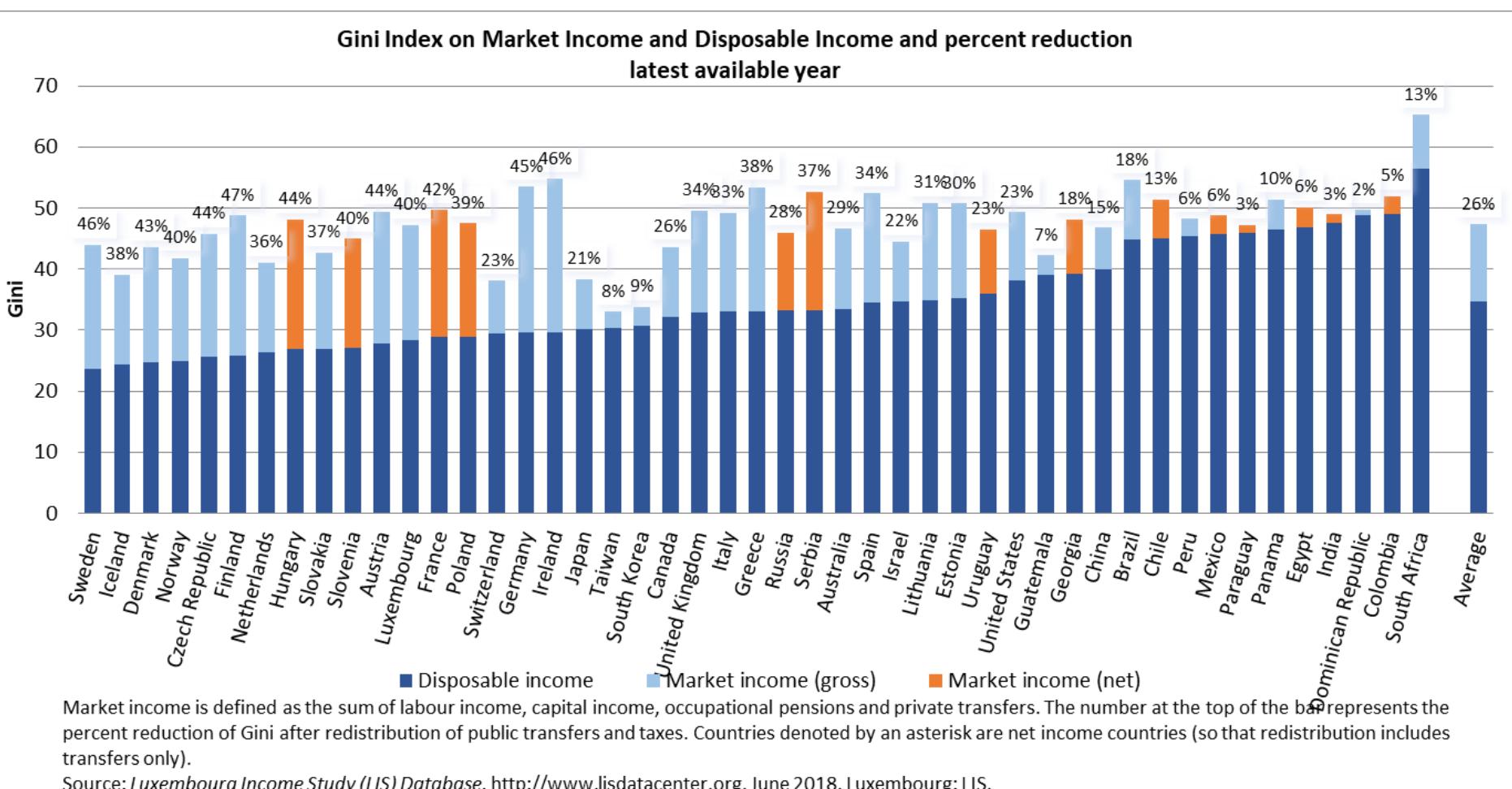
% of individuals with household income below 50% of median, latest available year



Source: LIS Inequality and Poverty Key Figures, <http://www.lisdatacenter.org>. June 2018. Luxembourg: LIS.

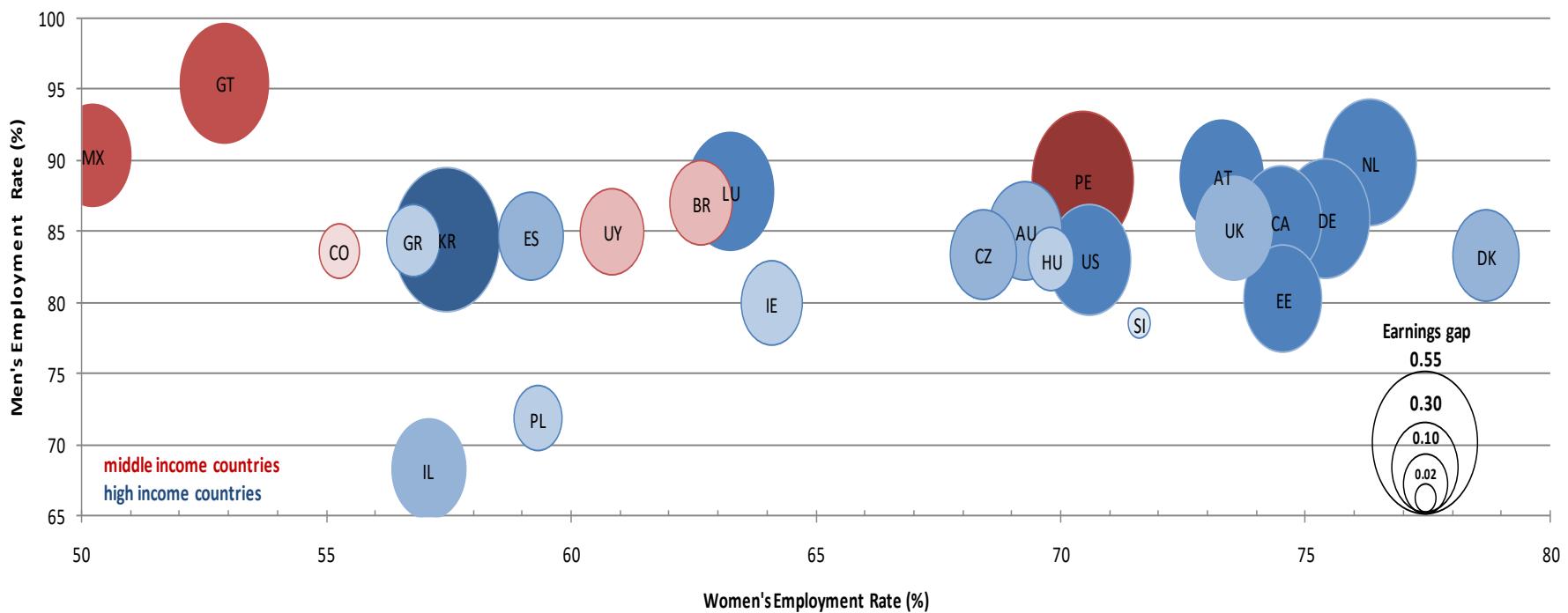
# Researching Policy Impacts

## Inequality reducing effect of Redistribution



# Comparing Employment Outcomes

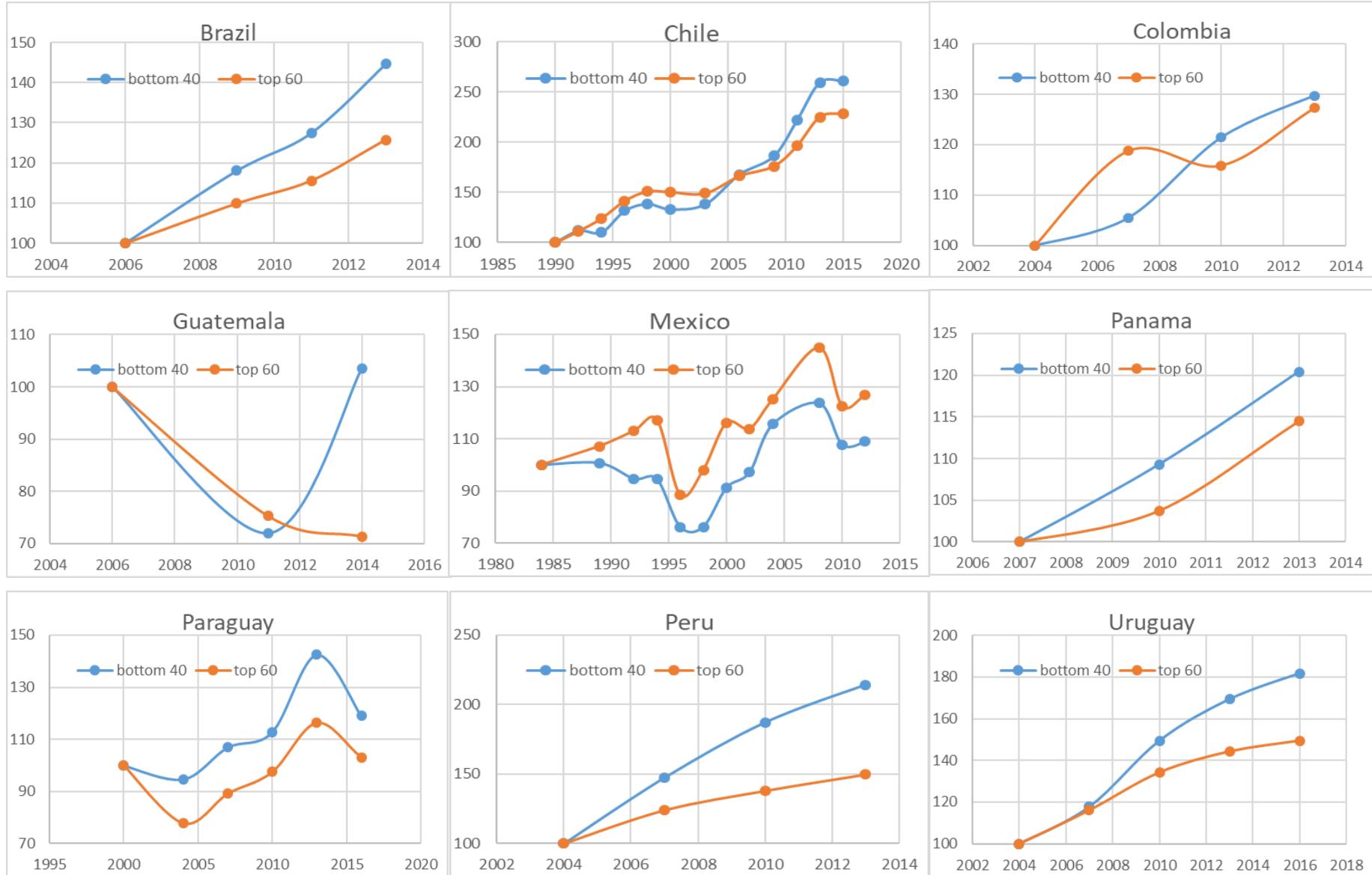
## Earnings Equality between Women and Men



L I S

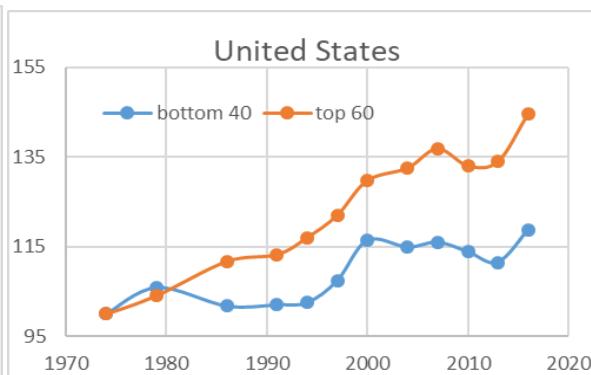
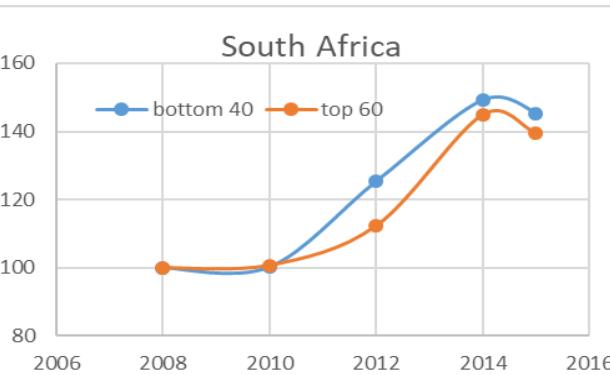
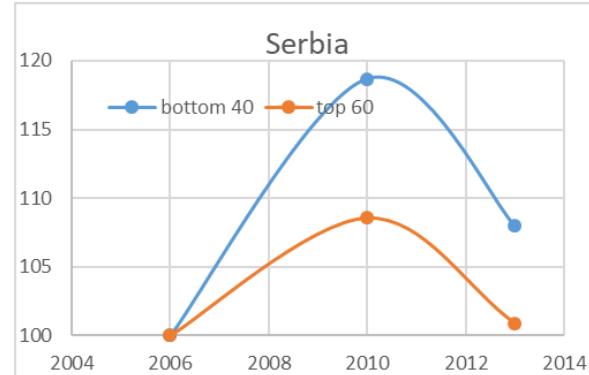
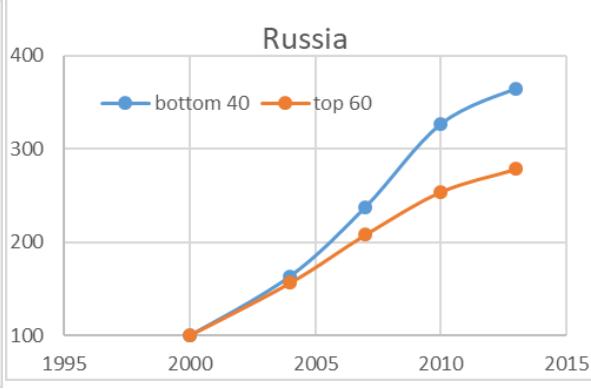
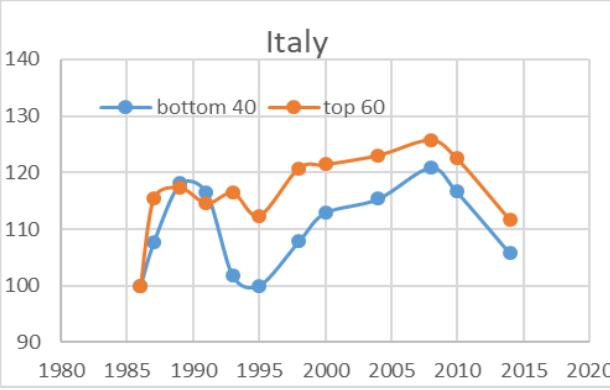
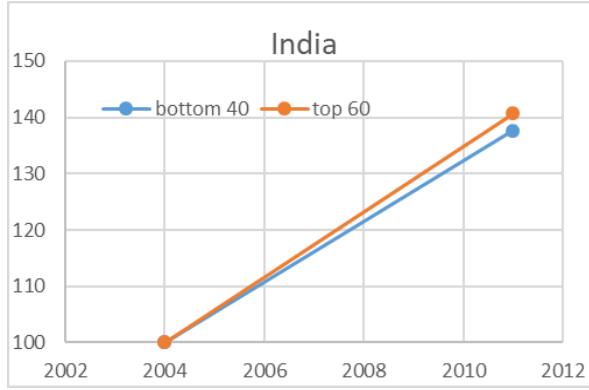
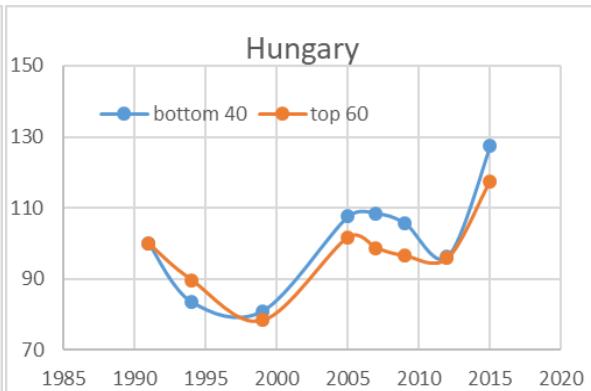
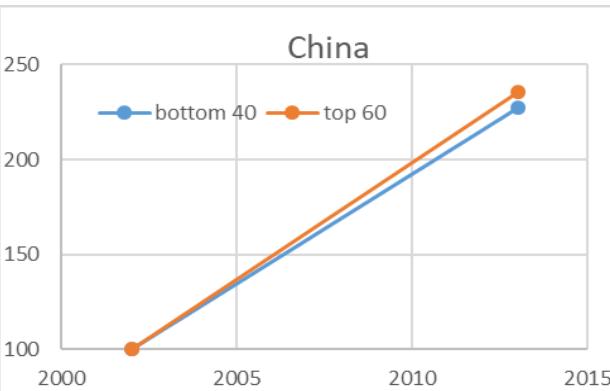
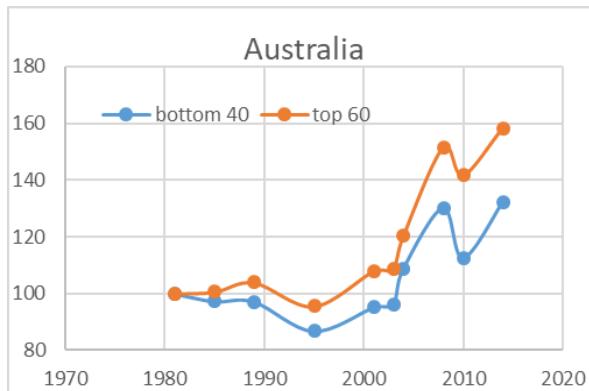
# Measuring SDGs: goal 10.1

By 2030, progressively achieve and sustain income growth of the bottom 40 per cent of the population at a rate higher than the national average



# Measuring SDGs: goal 10.1

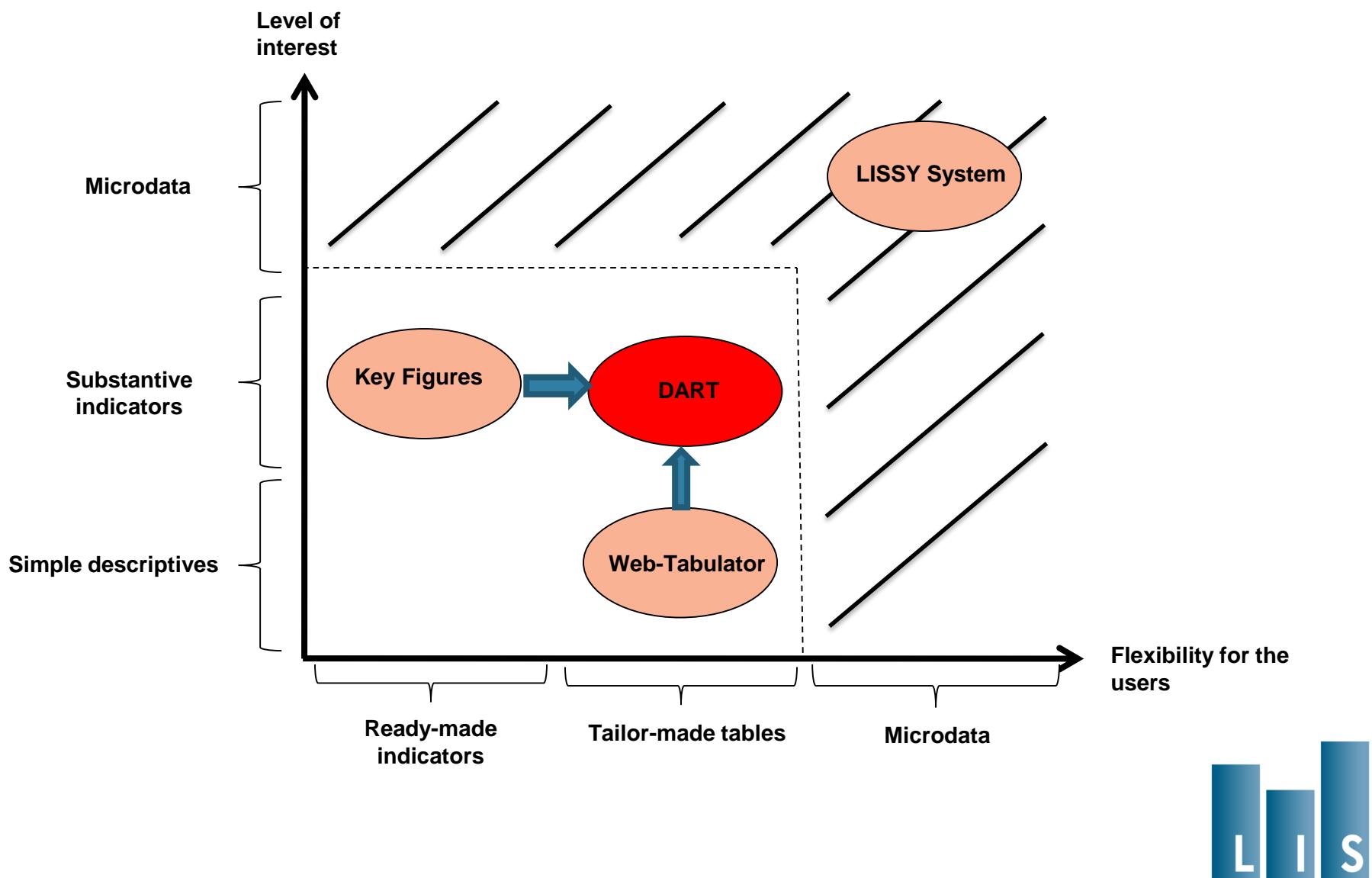
By 2030, progressively achieve and sustain income growth of the bottom 40 per cent of the population at a rate higher than the national average



# Part IV : Data dissemination & Documentation

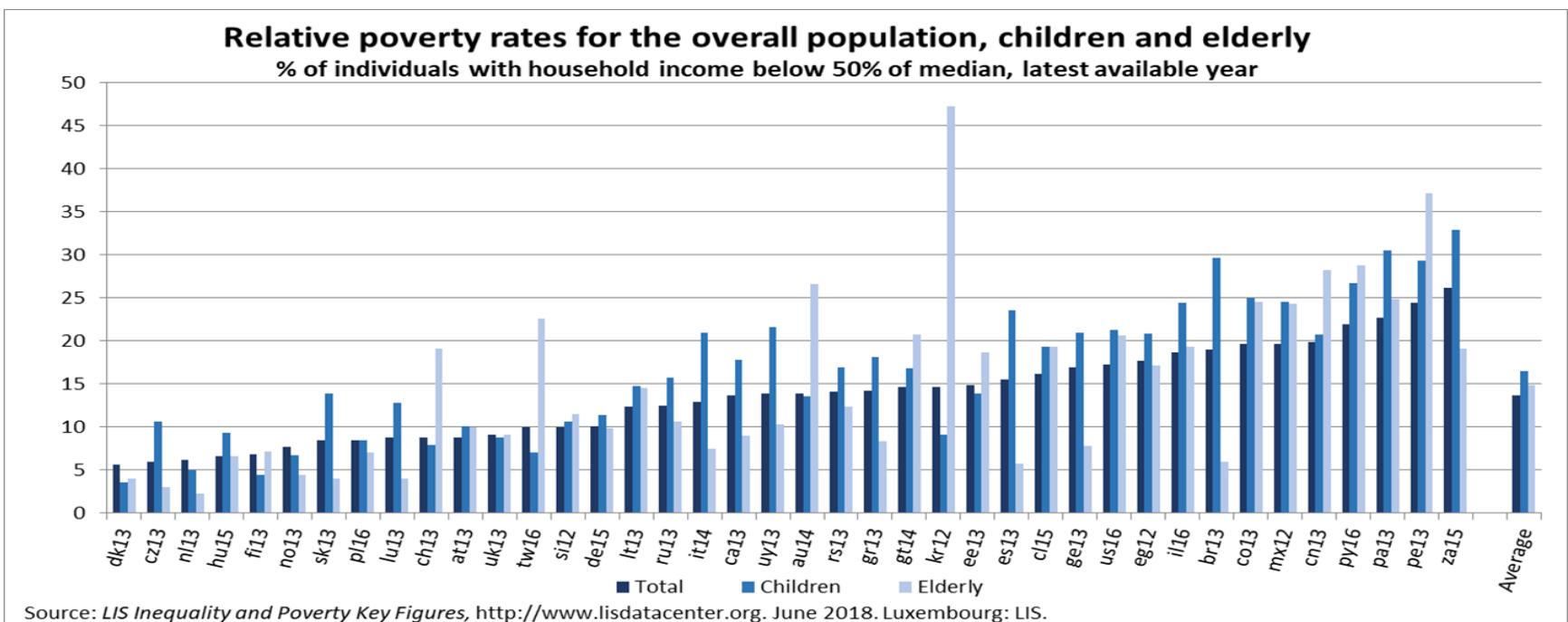
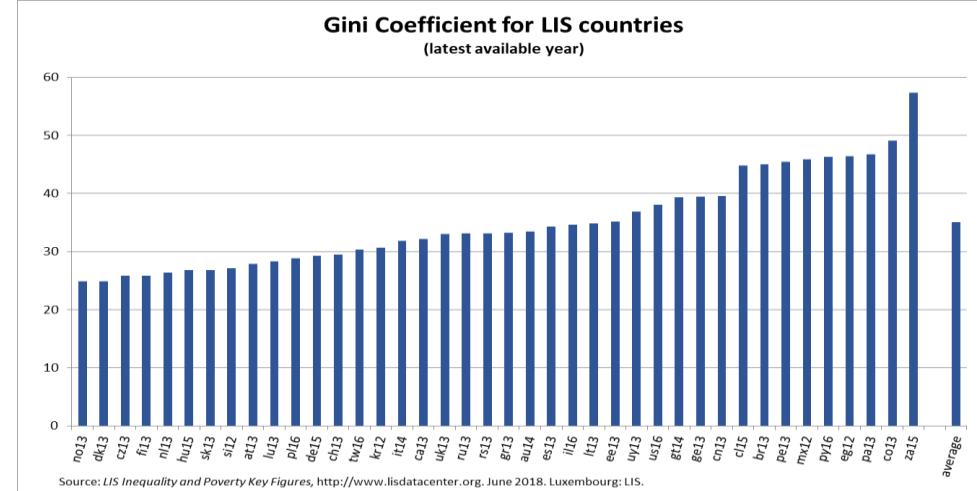


# Data access



# LIS Key Figures

- Multiple country-level **inequality** measures (e.g., Gini and Atkinson coefficients, percentile ratios)
- Relative **poverty** rates for various demographic groups
- Median and mean disposable household income



# LISSY: remote-execution system

- Fully automated, running 24 hours/day and 7 days/week
- Job Submission Interface (**JSI**): send statistical batch programs (SAS, SPSS, Stata or R) automatically processed and reports back aggregated results
- Micro-databases **cannot be downloaded** and no direct access to the data is permitted. **Only aggregated results** from statistical requests are sent back to the users.
- Access to LISSY is granted to researchers, incl. students, working for an academic, government or non-profit organization under the condition that use of the micro-data is restricted to **research purposes** only



# How LISSY works 1/2

The LISSY website homepage features a large group photo of researchers. Below the photo are three main sections: 'Login to LISSY' (with a lock icon), 'Login to Web Tabulator' (with a lock icon), and 'Key Figures' (with a bar chart icon). Each section has a brief description and a 'register' or 'login' button.

The LISSY Web User Interface shows a job submission form. The 'project' dropdown is set to 'LIS'. The 'statistical package' dropdown is set to 'R'. The 'subject' field contains 'SELF-TEACHING PACKAGE R PART II - REVIEW - EX8'. The right side of the interface shows a table with columns: job, project, processor, date, subject, and status. Below the form is a block of R code:

```
get_stack <- function(datasets, varp, varh, subset) {
  # READ DATASETS
  pp <- read.LIS(paste(datasets, 'p', sep = ""), labels = FALSE, vars = varp, subset = subset)
  hh <- read.LIS(paste(datasets, 'h', sep = ""), labels = FALSE, vars = varh)
  df <- merge(pp, hh, by = c("dname", "hid"))
  # MAP NEW VARIABLES
  df$homeowner <- ifelse(df$own %in% 100:199, 1, ifelse(df$own %in% 200:299, 0, NA))
  df$child <- ifelse(df$children %in% c(140, 200) , 0, ifelse(df$children == 110 , 1, ifelse(is.na(df$children), NA, 2)))
  df$ychild <- ifelse(df$child == 1 , 1, ifelse(is.na(df$child) , NA, 0))
  df$ochild <- ifelse(df$child == 2 , 1, ifelse(is.na(df$child) , NA, 0))
  df$partner <- ifelse(df$partner %in% c(100, 120, 200), 0, ifelse(is.na(df$partner) , NA, 1))
  df$meduc <- heduc <- 0
  df$meduc <- ifelse(df$educ == 2 , 1, ifelse(is.na(df$educ) , NA, 0))
  df$heduc <- ifelse(df$educ == 3 , 1, ifelse(is.na(df$educ) , NA, 0))
  df$ppp <- ifelse(df$dname == 'de04', 0.74, ifelse(df$dname == 'gr04', 0.62, 1))
  df$hrgw <- df$gross1
  df$hrgw <- ifelse(df$hrgw <= 0, NA, df$hrgw)
  df$germany <- ifelse (df$dname == 'de04', 1, 0)
  df$greece <- ifelse (df$dname == 'gr04', 1, 0)
  for (i in 1:length(datasets)) {
    df$hrgw <- df$hrgw / df$ppp
    topline <- with(df[lis.na(df$hrgw) & df$dname==datasets[i,], 10 * wNtile(hrgw,ppopwgt,0.5))
    df$hrgw <- with(df[lis.na(df$hrgw) & df$dname==datasets[i,], ifelse(df$hrgw>topline,topline,df$hrgw))
  }
  return(df)
}
```



# How LISSY works 2/2

LISSY WEB USERINTERFACE

submit job

job session today jobs job library

refresh from 19 Apr 18 - 19 Apr 18 advanced search

view job discard job

ge	date	subject
	19 Apr 18 17:05	SELF-TEACHING PACKAGE R PART II - REVIEW - EX7
	19 Apr 18 16:43	SELF-TEACHING PACKAGE R PART II - REVIEW - EX8
	19 Apr 18 16:40	SELF-TEACHING PACKAGE R PART II - REVIEW - EX1
	19 Apr 18 15:45	SELF-TEACHING PACKAGE Stata PART II - REVIEW - EX8
	19 Apr 18 14:11	SELF-TEACHING PACKAGE Stata PART II - REVIEW - EX2
	19 Apr 18 09:35	SELF-TEACHING PACKAGE Stata PART II - REVIEW - EX7
	19 Apr 18 09:33	SELF-TEACHING PACKAGE Stata PART II - REVIEW - EX2
	19 Apr 18 09:27	SELF-TEACHING PACKAGE Stata PART II - REVIEW - EX2
	19 Apr 18 09:10	SELF-TEACHING PACKAGE SPSS PART II - REVIEW - EX7
	19 Apr 18 09:02	SELF-TEACHING PACKAGE SPSS PART II - REVIEW - EX6
	19 Apr 18 08:42	SELF-TEACHING PACKAGE SPSS PART II - REVIEW - EX5
	19 Apr 18 08:36	SELF-TEACHING PACKAGE SPSS PART II - REVIEW - EX4
	19 Apr 18 08:24	SELF-TEACHING PACKAGE SPSS PART II - REVIEW - EX3

19 Apr 18 16:43 SELF-TEACHING PACKAGE R PART II - REVIEW - EX8

job text listing

Call:  
glm(formula = model, data = df, weights = df\$pwgt, subset = sex ==  
s)

Deviance Residuals:

Min	1Q	Median	3Q	Max
-3.05054	-0.06792	-0.00059	0.06947	2.08506

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )		
(Intercept)	1.318e+00	7.205e-02	18.295	< 2e-16 ***		
age	4.978e-02	3.710e-03	13.417	< 2e-16 ***		
I(age^2)	-4.973e-04	4.604e-05	-10.803	< 2e-16 ***		
meduc	1.844e-01	8.555e-03	21.556	< 2e-16 ***		
heduc	5.242e-01	8.911e-03	58.827	< 2e-16 ***		
immigr	-1.169e-01	7.877e-03	-14.855	< 2e-16 ***		
partn	3.108e-02	7.947e-03	3.911	9.22e-05 ***		
ychild	7.717e-02	8.188e-03	9.424	< 2e-16 ***		
ochild	6.744e-02	7.319e-03	9.215	< 2e-16 ***		
ptime	-3.008e-01	1.469e-02	-20.477	< 2e-16 ***		
homeowner	1.371e-01	6.210e-03	22.078	< 2e-16 ***		
germany	7.043e-01	6.535e-03	107.765	< 2e-16 ***		
greece	5.909e-01	7.206e-03	82.001	< 2e-16 ***		
---						
Signif. codes:	0 ****	0.001 ***	0.01 **	0.05 *	0.1 .	1
(Dispersion parameter for gaussian family taken to be 0.02541886)						
Null deviance: 1420.19 on 34777 degrees of freedom						
Residual deviance: 883.69 on 34765 degrees of freedom						
(10071 observations deleted due to missingness)						
AIC: 76786						
Number of Fisher Scoring iterations: 2						

define search parameters X

(if using multiple keywords, separate them by a ' ')

subject line contains

job contains

search conditions  are casesensitive

scope  include jobs that have been discarded

search cancel

# METADATA INFORMATION SYSTEM (METIS)

**METIS enables browsing, aggregating, and exporting LIS/LWS Databases documentation tailored to the users' needs**

## 1. Overview of datasets and variables

Overview of the contents of the LIS/LWS databases in terms of datasets and variables

## 2. Compare datasets

Select datasets and compare generic information among them (original survey, Social Security, Key Figures)

## 3. Compare variables

Select variables and compare generic information among them (variables definitions and standard labels)

## 4. Cross-compare – main functionality

View the availability of the selected variables in the selected datasets and compare dataset-specific information (statistics and notes)



## SELECTION RESULTS

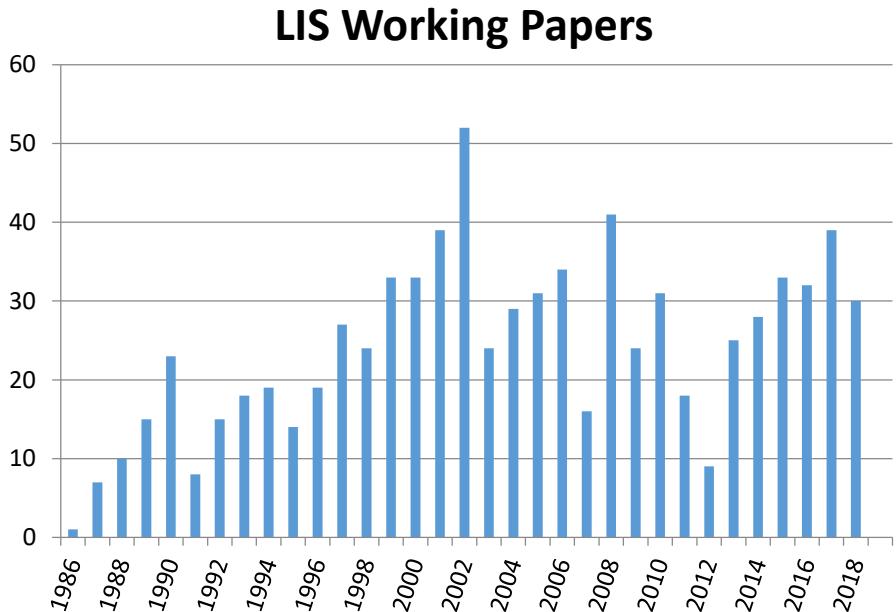
## Select datasets      Select variables

# LIS/LWS Working papers

## Working Paper Series

- 813 WP of which over 400 published in peer reviewed journals
- Aldi Award: best paper by an under 40 author
- included on RePec (Research Papers in Economics)

- **New WP:** 69 new WP added in 2017 and 2018.
- **RePEC:** Since May 2017, LIS/ LWS working papers have been downloaded more than 9,000 times
- **H-index:** LIS = 166  
LWS = 50



# How do we respond to those challenges?



# V. LIS golden rules for harmonisation

- Set clear definitions for LIS variables
  - *Maximise comparability by setting clear definitions for each variable (and trying to stick to them as much as possible)*
  - *Document very well any deviation from the general definition*
- Complement easiness of use with flexibility of use
  - *Enhance user-friendliness by providing fully standardised variables (standard variables, recodes, dummies, aggregate variables)*
  - *Allow users the flexibility to create other concepts by leaving a large amount of detailed information*
- Adapt the LIS template to the changing environment (over time and space)
  - *Template revisions*
  - *Revisions of previous datasets*



Overall guiding principle: **OPERATIONAL COMPARABILITY**



# Final remarks

- Sample coverage
- Individual level income (high non response affecting overall household income).
- Constructing upper level aggregated variables to help in the consistency checks
- Clearer documentation and definitions