

**nfdi 4 health** funded by **DFG** Deutsche Forschungsgemeinschaft German Research Foundation

**STRATOS INITIATIVE** TG3

**Universitätsmedizin GREIFSWALD**

## Accessible structured Initial Data Analysis (IDA) and Data Quality Assessments (DQA)

Carsten Oliver Schmidt, Lara Lusa,  
Marianne Huebner on behalf of  
STRATOS TG3

**DAGSTat 2025**

Heubner et al. BMC Medical Research Methodology (2018) 18:104  
https://doi.org/10.1186/s13227-018-0593-y

BMC Medical Research Methodology

**RESEARCH ARTICLE** Open Access

Hidden analyses: a review of reporting practice and recommendations for more transparent reporting of initial data analyses

Marianne Huebner<sup>1,2\*</sup>, Werner Vach<sup>3</sup>, Saskia le Cessie<sup>4</sup>, Carsten Oliver Schmidt<sup>5</sup>, Lara Lusa<sup>6,7</sup> and on behalf of the Topic Group "Initial Data Analysis" of the STRATOS Initiative (STRengthening Analytical Thinking for Observational Studies). <http://www.stratos-initiativestudies.com>

**Research Square**

Observational Studies 4 (2018) 171–192 Submitted 7/17; Published 4/18

**A Contemporary Conceptual Framework for Initial Data Analysis**

Marianne Huebner [huebner@stt.msu.edu](mailto:huebner@stt.msu.edu)  
Department of Statistics and Probability  
Michigan State University  
East Lansing, MI 48824, USA

Saskia le Cessie [s.lecessie@lumc.nl](mailto:s.lecessie@lumc.nl)  
Department of Clinical Epidemiology and Department of Medical Statistics and Bioinformatics  
Leiden University Medical Center  
Leiden, The Netherlands

**Regression without regrets – initial data analysis is an essential prerequisite to multivariable regression**

Georg Heinze ([georg.heinze@meduniwien.ac.at](mailto:georg.heinze@meduniwien.ac.at))  
Medical University of Vienna.  
Mark Baillie  
Novartis (Switzerland)  
Lara Lusa  
University of Primorska  
Willi Sauerbrei  
University of Freiburg  
Carsten Oliver Schmidt

**PLOS ONE**

**STRATOS TG3**

**PLUS COMPUTATIONAL BIOLOGY**

**RESEARCH ARTICLE**  
Initial data analysis for longitudinal studies to build a solid foundation for reproducible analysis

Lara Lusa<sup>1,2,\*</sup>, Cécile Proust-Lima<sup>3</sup>, Carsten O. Schmidt<sup>4</sup>, Katherine J. Lee<sup>5,6</sup>, Saskia le Cessie<sup>4</sup>, Mark Baillie<sup>7</sup>, Frank Lawrence<sup>8,9</sup>, Marianne Huebner<sup>10,11</sup>, on behalf of TG3 of the STRATOS Initiative<sup>9</sup>

**EDITORIAL**  
**Ten simple rules for initial data analysis**

Mark Baillie<sup>1</sup>, Saskia le Cessie<sup>2</sup>, Carsten Oliver Schmidt<sup>3</sup>, Lara Lusa<sup>4</sup>,  
Marianne Huebner<sup>10,11\*</sup>, for the Topic Group "Initial Data Analysis" of the STRATOS Initiative<sup>9</sup>

<sup>1</sup> Novartis, Basel, Switzerland, <sup>2</sup> Department of Clinical Epidemiology and Department of Biomedicine, University of Münster, Münster, Germany, <sup>3</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, <sup>4</sup> Department of Biostatistics, Leiden University Medical Center, Leiden, The Netherlands, <sup>5</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, <sup>6</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, <sup>7</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, <sup>8</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, <sup>9</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, <sup>10</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, <sup>11</sup> Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America

**1** Novartis, Basel, Switzerland, **2** Department of Clinical Epidemiology and Department of Biomedicine, University of Münster, Münster, Germany, **3** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, **4** Department of Biostatistics, Leiden University Medical Center, Leiden, The Netherlands, **5** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, **6** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, **7** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, **8** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, **9** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, **10** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America, **11** Department of Biostatistics, University of Michigan, Ann Arbor, MI, United States of America

## Aim of IDA



The aim of IDA is to provide a data set and reliable findings on this data set which allows researchers to work with this data set in a responsible manner.

Huebner et al. 2018

## Common (implicit) assumptions underlying analyses

1. Informative data
2. Sufficient data quantity
3. Limited missing data
4. Little systematic noise
5. Correct data labels
6. Class balance
7. Non-adversarial data
8. Independent and Identically Distributed (i.i.d.)
9. Linearity
10. Homoscedasticity
11. Non-multicollinearity
12. Consistent scaling of features
13. Sample representativeness

.....  
Background information  
Study design  
Variable descriptions  
Range violations  
Contradictions  
Inadmissible values  
Volume  
Unit / item missingness  
Missing patterns  
Missing mechanisms  
Univariate descriptions  
Multivariate descriptions  
Associations  
.....

## Data Quality Framework for EU medicines regulation 2023



30 October 2023  
Data Analytics and Methods Task Force  
EMA/326985/2023

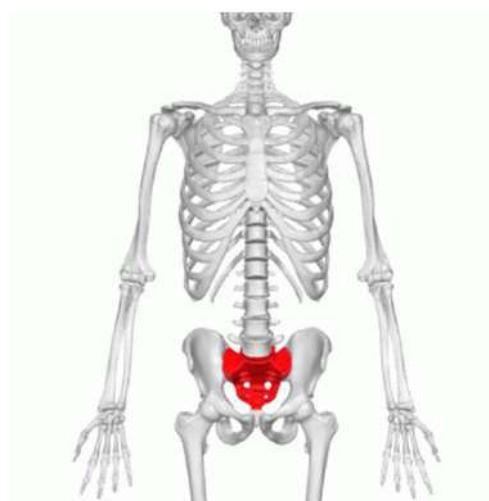
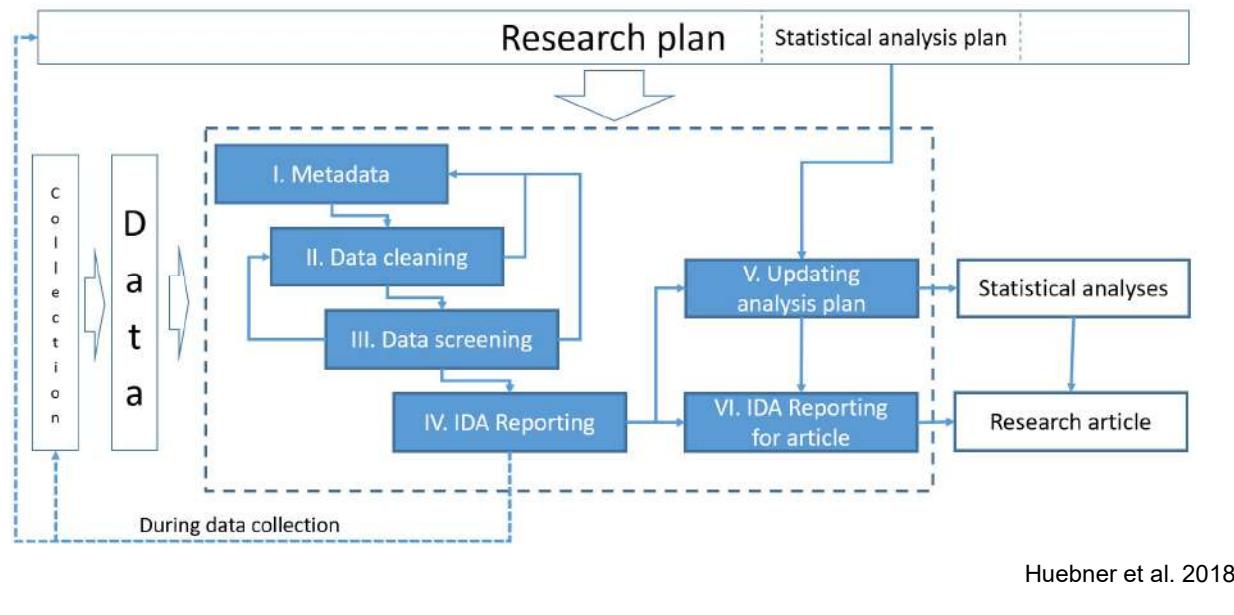


## What does accessible IDA and DQA mean?



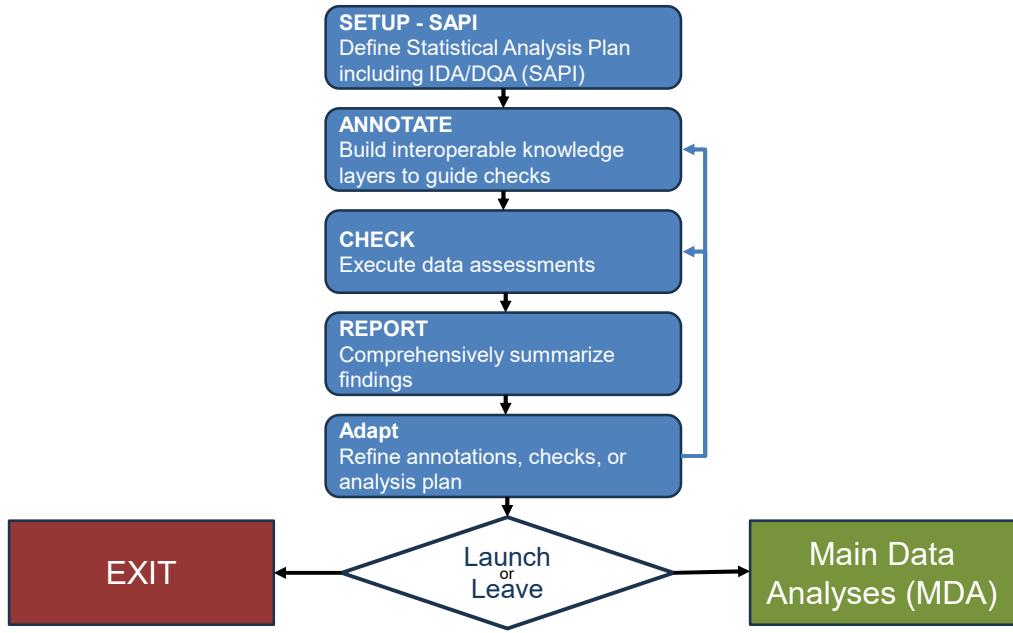
A transparent approach to planning, checking, reporting, data handling, and decision-making processes that govern the creation and use of analysis datasets as part of the Main Data Analysis (MDA).

## Initial Data Analysis Workflow

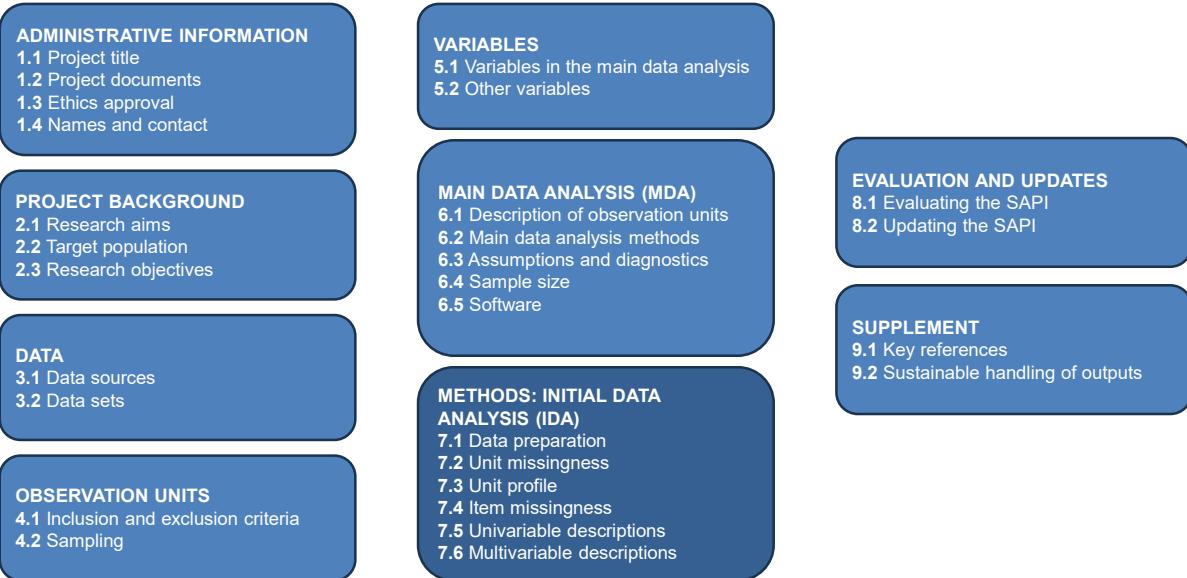


<https://en.wikipedia.org/wiki/Sacrum>

## The SACRAL-Cycle



## Setup – Define Statistical Analysis Plan including IDA/DQA (SAPI)



## Setup – Define Statistical Analysis Plan including IDA/DQA (SAPI)

Table 1. Initial data analysis checklist for data screening in longitudinal studies.

Topic	Item	Features
<b>IDA screening domain: Participation profile</b>		
Time frame	P1	Provide number of time points and intervals at which measurements are taken, using the time metric that best reflects the time from inclusion in the study, or calendar time in studies that involve long enrollment times. Highlight the differences between the time of first measurements and follow-up times.
Time metric	P2	Describe the time metric and corresponding time points specified in the analysis strategy, if different from the time metric described in P1.
Participants	P3	Provide the number of participants who attended the assessment by time metric(s).
<b>Optional extensions: Participation Profile</b>		
Other time metrics	PE1	Use different time metric(s) to describe the time frame of the study, if applicable and appropriate, e.g. calendar time or data collection visits.
<b>IDA screening domain: Missing data (outcome variable and explanatory variables)</b>		
Non-enrollment	M1	Describe the non-enrolled, i.e., the participants that were selected but did not enter the study (and the reasons, if available), if applicable.
Drop-out	M2	Describe the participants who dropped out from the study during the follow-up (loss to follow-up and other possible reasons: death, withdrawal, missing by design, if applicable).
Intermittent visit missingness	M3	Describe the participants that have missing data for some of the measurements (intermittent, occasional omission, but do not drop out out of the study).

Lusa et al. 2024

## Setup – Define Statistical Analysis Plan including IDA/DQA (SAPI)

Data sets      3.2

### Describe how the data is provided for this analysis project, for example format and content of data sets

- Describe how the data was processed from their raw state after data collection until transfer for use in this analysis project. For example, any preliminary data editing, deletion of cases and variables, or rules underlying the computation of new variables.

## Annotate – Build interoperable knowledge layers to guide checks



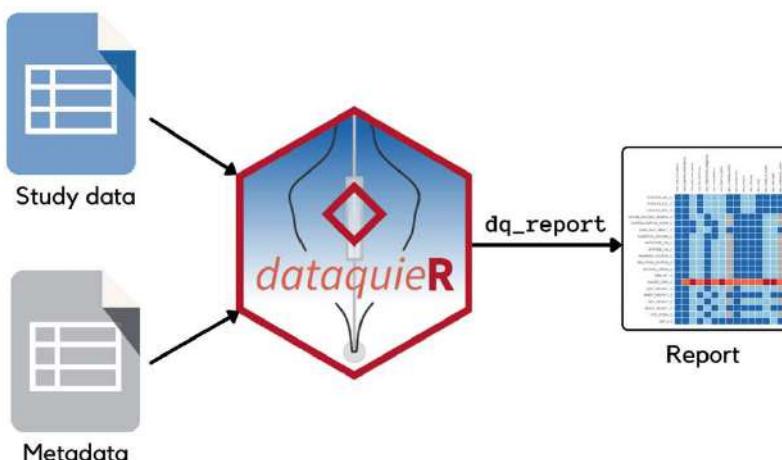
Study data



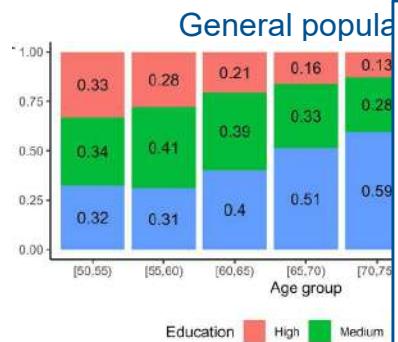
Metadata

A	B	C	D	E	F	G	
VAR_NAMES	LABEL	DATA_TYPE	SCALE_LEVEL	VALUE_LABELS	MISSING_LIST_TABLE	HARD_LIMITS	DETECT
v00000	CENTER_0	integer	nominal	1 - Berlin   2 - Hamburg   3 - Leipzig   4 - Cologne   5 - Munich			
v00001	PSEUDO_ID	string	na				
v00002	SEX_0	integer	nominal	0 - females   1 - males		[18,Inf)	
v00003	AGE_0	integer	ratio				
v00103	AGE_GROUP_0	string	ordinal				
v01003	AGE_1	Integer	ratio			[18,Inf)	
v01002	SEX_1	integer	nominal	0 = females   1 = males			
v10000	PART_STUDY	integer	nominal	0 = no   1 = yes			
v00004	SBP_0	float	ratio		missing_table	[80;180]	[0;265]
v00005	DBP_0	float	ratio		missing_table	[50;Inf)	[0;265]
v00006	GLOBAL_HEALTH_V	float	ratio		missing_table	[0;10]	
v00007	ASTHMA_0	integer	nominal	0 = no   1 = yes	missing_table		[0;1]
v00008	VO2_CAPCAT_0	string	ordinal	A = excellent < B = good	missing_table		
v00009	ARM_CIRC_0	float	ratio		missing_table		
v00109	ARM_CIRC_DISC_0	integer	ordinal	1 = (-Inf,20] < 2 = (20,30]	missing_table		[1;3]
v00010	ARM_CUFF_0	integer	ordinal	1 = (-Inf,20] < 2 = (20,30]	missing_table		[1;3]
v00011	USR_VO2_0	string	nominal	USR_101   USR_103   U	missing_table		
v00012	USR_BP_0	string	nominal	USR_121   USR_123   U	missing_table		
v00013	EXAM_DT_0	datetime	interval				[2018-01-01 00:00:00 CET]
v20000	PART_PHYS_EXAM	integer	nominal	0 = no   1 = yes			
v00014	CRP_0	float	ratio		missing_table	[0;Inf)	[0;16;Inf]
v00015	BSG_0	float	ratio		missing_table	[0;100]	
v00016	DEV_NO_0	integer	nominal				
v00017	LAB_DT_0	datetime	interval				[2018-01-01 00:00:00 CET]
v20000	DATA_LAB	integer	nominal	0 = no   1 = yes			

## Check – Execute data assessments



## Report – Comprehensively summarize findings



Lusa et al. 2024

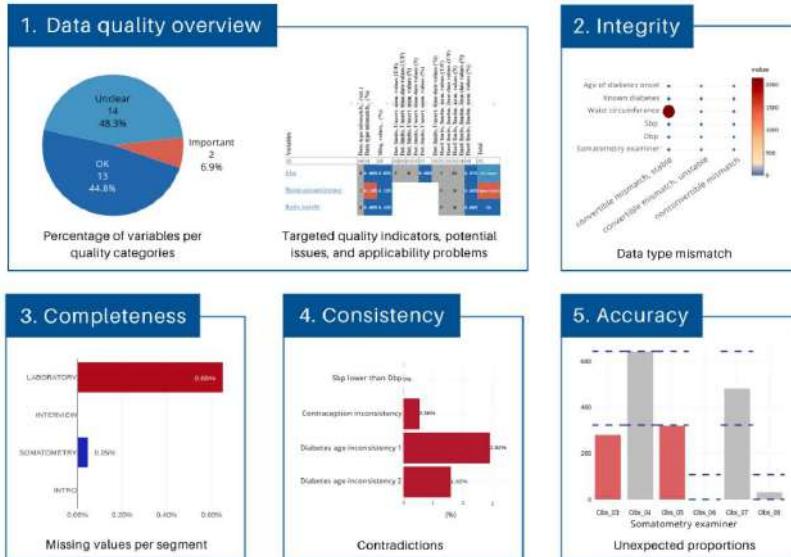
1 Preface
2 Initial data analysis and data screening checklist for longitudinal data
3 SHARE data description
4 Data example (from research question to IDA plan)
5 Data screening
5.1 Participation profile
5.2 Missing values
5.2.1 Non-enrollment (M1)
5.2.2 Drop-out (M2) and intermittent missingness (M3)
5.2.3 Variable missingness (item missingness, M4)
5.2.4 Patterns (M5)
5.2.5 Comparison of non-enrolled and target population (ME1)
5.2.6 Probability of loss to follow-up and death (ME2)
5.2.7 Dropout effect on outcome (ME3)
5.3 Univariate descriptions
5.4 Multivariate description of data
5.5 Longitudinal aspects

Baseline characteristics by type of missingness.

	N	Complete N=2681	Death N=978	Intermittent missing N=476
gender :				
Female	5452	0.54 <sup>140/2681</sup>	0.51 <sup>49/978</sup>	0.50 <sup>242/476</sup>
age_int	5452	52.00 58.00 66.00	68.00 75.00 81.00	52.00 58.00 64.00
60-69		60.28 ± 8.79	73.29 ± 10.44	59.55 ± 8.18
70-80				
80+				
age_int_cat :	5452	0.54 <sup>140/2681</sup>	0.12 <sup>12/978</sup>	0.59 <sup>38/476</sup>
60-69		0.29 <sup>78/2681</sup>	0.21 <sup>11/978</sup>	0.27 <sup>12/476</sup>
70-80		0.14 <sup>34/2681</sup>	0.41 <sup>18/978</sup>	0.13 <sup>5/476</sup>
80+		0.02 <sup>6/2681</sup>	0.26 <sup>5/978</sup>	0.01 <sup>1/476</sup>
weight	5361	66.0 76.0 86.0	62.5 71.0 81.0	65.0 76.0 85.0
77.2 ± 15.2		72.7 ± 15.0	77.1 ± 15.6	
height_imp	5418	178.00	175.00	178.00
171.82 ± 9.04		169.34 ± 8.80	171.66 ± 8.98	
education_imp	5428	0.17 <sup>44/2681</sup>	0.38 <sup>11/978</sup>	0.19 <sup>9/476</sup>
Low		0.38 <sup>10/2681</sup>	0.39 <sup>11/978</sup>	0.41 <sup>10/476</sup>
Medium		0.45 <sup>12/2681</sup>	0.23 <sup>5/978</sup>	0.40 <sup>12/476</sup>
High				
pa_vig_freq	5423	0.67 <sup>118/2681</sup>	0.35 <sup>33/978</sup>	0.66 <sup>31/476</sup>
pa_low_freq	5422	0.94 <sup>201/2681</sup>	0.73 <sup>73/978</sup>	0.95 <sup>46/476</sup>
cusmoke_imp	5423	0.22 <sup>58/2681</sup>	0.34 <sup>12/978</sup>	0.27 <sup>12/476</sup>
Yes				
maxgrip	5272	29.0 36.0 48.0	21.5 29.0 38.0	28.0 37.5 49.0
38.5 ± 12.5		30.3 ± 11.9	38.5 ± 13.1	

a b c represent the lower quartile a, the median b, and the upper quartile c for continuous is the number of non-missing values.

## Report – Comprehensively summarize findings



<https://dataquality.qihs.uni-greifswald.de/>

## Adapt – Refine annotations, checks, or analysis plan

### ADMINISTRATIVE INFORMATION

- 1.1 Project title
- 1.2 Project documents
- 1.3 Ethics approval
- 1.4 Names and contact

### PROJECT BACKGROUND

- 2.1 Research aims
- 2.2 Target population
- 2.3 Research objectives

### DATA

- 3.1 Data sources
- 3.2 Data sets

### OBSERVATION UNITS

- 4.1 Inclusion and exclusion criteria
- 4.2 Sampling

### VARIABLES

- 5.1 Variables in the main data analysis
- 5.2 Other variables

### MAIN DATA ANALYSIS (MDA)

- 6.1 Description of observation units
- 6.2 Main data analysis methods
- 6.3 Assumptions and diagnostics
- 6.4 Sample size
- 6.5 Software

### METHODS: INITIAL DATA ANALYSIS (IDA)

- 7.1 Data preparation
- 7.2 Unit missingness
- 7.3 Unit profile
- 7.4 Item missingness
- 7.5 Univariable descriptions
- 7.6 Multivariable descriptions

### EVALUATION AND UPDATES

- 8.1 Evaluating the SAPI
- 8.2 Updating the SAPI

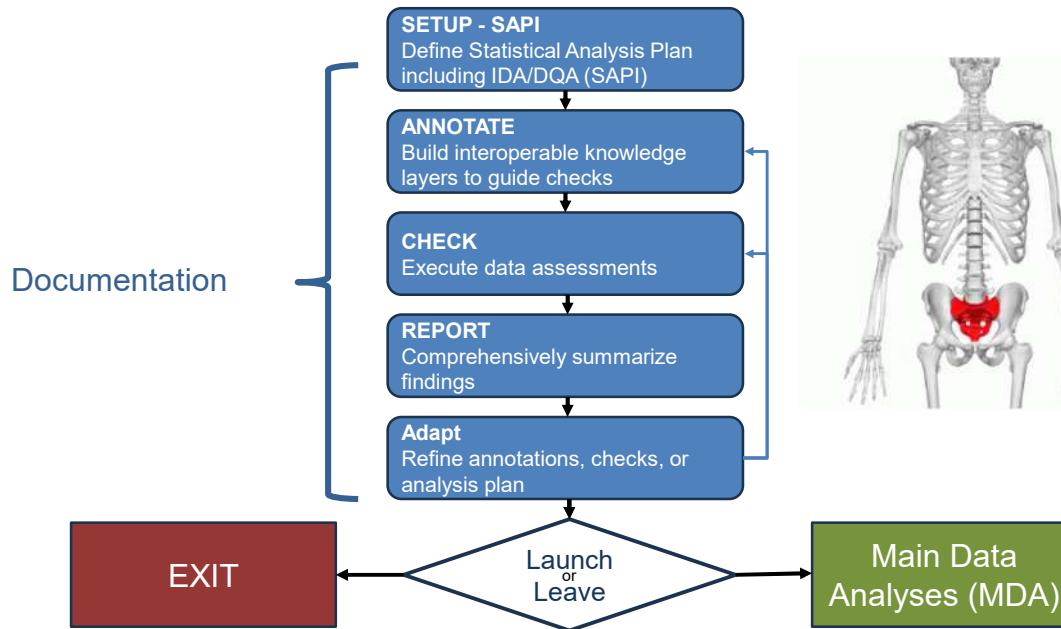
### SUPPLEMENT

- 9.1 Key references
- 9.2 Sustainable handling of outputs

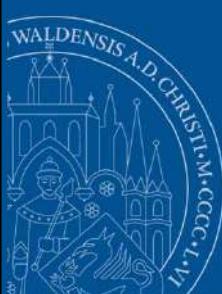
## Launch or Leave – Conduct or discard intended analyses

Section 63 Spatial and spatio-temporal Statistics I  02:00 pm - 03:20 pm Topic : 63 Spatial and Spatio-temporal Statistics Form : Oral Chair(s): Dr. Almond	Section 46 Advances in Item Response Theory-based measurement of abilities  02:00 pm - 03:20 pm Topic : 46 Testing and Scaling Form : Oral Chair(s): Prof. Dr. Andersen	Section 21 Tree-based AI  02:00 pm - 03:20 pm Topic : 21 Artificial Intelligence and Machine Learning Form : Oral Chair(s): Prof. Dr. Gähler	Section 11 Causal inference – mixed topics  02:00 pm - 03:20 pm Topic : 11 Causal Inference Form : Oral Chair(s): Dr. Claudia Böhmhofer & Institute Institute	Section 45 Robust Statistical Methods for Model Fitting and Data Analysis I  02:00 pm - 03:20 pm Topic : 45 Robust and Nonparametric Statistics Form : Oral
Break	Break	Break	Break	Break
Section 63 Spatial and spatio-temporal Statistics II  03:50 pm - 05:10 pm Topic : 63 Spatial and Spatio-temporal Statistics Form : Oral Chair(s): Dr. Almond	Section 46 Linking, equating, and norming in the context of IRT  03:50 pm - 05:10 pm Topic : 46 Testing and Scaling Form : Oral Chair(s): Prof. Dr. Andersen	Section 21 xAI, Interpretability  03:50 pm - 05:10 pm Topic : 21 Artificial Intelligence and Machine Learning Form : Oral Chair(s): Mariele Stute	Emil Gumbel: The prediction of extreme events (Documentary Movie)  03:50 pm - 05:10 pm Topic : 98 Special Session Form : Oral Chair(s): Mariele Stute	Section 45 Robust Statistical Methods for Model Fitting and Data Analysis II  03:50 pm - 05:10 pm Topic : 45 Robust and Nonparametric Statistics Form : Oral
Break	Break	Break	Break	Break

## Accessible.... Transparent..... Reproducible



Thanks to the many involved



**STRATOS**  
INITIATIVE



Carsten.schmidt@uni-greifswald.de  
Universitätsmedizin Greifswald . KöR  
ICM SHIP/KEF  
Fleischmannstraße 8 . 17475 Greifswald  
[www.medizin.uni-greifswald.de](http://www.medizin.uni-greifswald.de)

© Copyright 2024. All rights reserved.