

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email editorial.bmjopen@bmj.com

BMJ Open

Identifying and sharing data for secondary data analysis of physical activity, sedentary behaviour and their determinants across the life course in Europe: general principles and an example from DEDIPAC

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2017-017489
Article Type:	Research
Date Submitted by the Author:	01-May-2017
Complete List of Authors:	Lakerveld, Jeroen; VU University Medical Center, Loyen, A; VU University Medical Center, Ling, Fiona; University of Limerick, Physical Education and Sport Sciences Department DeCraemer, Marieke; Ghent University, movement and sport science van der Ploeg, Hidde; VU University Medical Center Amsterdam, Department of Public and Occupational Health; Prevention Research Collaboration, Sydney School of Public Health, University of Sydney O'Gorman, Donal; Dublin City University, Centre for Preventive Medicine Carlin, Angela; University of Limerick, Centre for physical activity and health research Caprinica, Laura; University of Rome Foro Italico Kalter, Joeri; VU University Medical Center Oppert, Jean-Michel; University Pierre et Marie Curie, Institute of cardiometabolism and nutrition (ICAN), Pitie-Salpêtrière hospital (AP-HP) Chastin, Sebastian; Glasgow Caledonian University, Institute of Applied Health Research, School of Health and Life Sciences; Ghent University, Department of Movement and Sports Sciences Cardon, Greet; Ghent University, Movement and Sports Sciences Brug, Johannes; VU University Medical Center, EMGO+, Institute for Health and Care Research MacDonncha, Ciaran; University of Limerick, Centre for Physical Activity and Health Research
Primary Subject Heading:	Epidemiology
Secondary Subject Heading:	Research methods, Public health
Keywords:	EPIDEMIOLGY, PREVENTIVE MEDICINE, PUBLIC HEALTH

SCHOLARONE™
Manuscripts

1
2
3 Identifying and sharing data for secondary data analysis of physical activity,
4 sedentary behaviour and their determinants across the life course in Europe:
5 general principles and an example from DEDIPAC
6
7

8 Jeroen Lakerveld^{1§}, Anne Loyen¹, Fiona Ling^{2,3}, Marieke De Craemer⁴, Hidde P. van der Ploeg⁵, Donal
9 J. O’Gorman⁶, Angela Carlin², Laura Capranica⁷, Joeri Kalter¹, Jean-Michel Oppert⁸, Sebastien
10 Chastin⁹, Greet Cardon⁴, Johannes Brug^{1,10}, Ciaran MacDonncha²
11
12

13
14
15 ¹ Department of Epidemiology and Biostatistics, EMGO Institute for Health and Care Research, VU
16 University medical center, De Boelelaan 1089a, 1081, HV, Amsterdam, The Netherlands.
17

18 ² Centre for Physical Activity and Health Research, Department of Physical Education and Sport
19 Sciences, University of Limerick, Ireland.
20

21 ³ Institute of Sport, Exercise and Active Living, Victoria University, Melbourne, Australia.
22

23 ⁴ Department of Movement and Sports Sciences, Ghent University, Ghent, Belgium.
24

25 ⁵ Department of Public and Occupational Health, EMGO Institute for Health and Care Research, VU
26 University medical center, Amsterdam, The Netherlands.
27

28 ⁶ Centre for Preventive Medicine, School of Health & Human Performance, Dublin City University,
29 Dublin 9, Republic of Ireland.
30

31 ⁷ Department of Movement, Human and Health Sciences, University of Rome Foro Italico, Rome,
32 Italy.
33

34 ⁸ Institute of Cardiometabolism and Nutrition (ICAN), University Pierre et Marie Curie-Paris6;
35 Department of Nutrition, Pitie-Salpetriere hospital (AP-HP), Paris, France.
36

37 ⁹ Institute of Applied Health Research, School of Health and Life Science, Glasgow Caledonian
38 University, Glasgow, UK.
39

40 ¹⁰ Amsterdam School for Communication Research, University of Amsterdam, Amsterdam, the
41 Netherlands.
42
43
44

45
46
47
48
49 **§ Corresponding author:** Jeroen Lakerveld; EMGO Institute for Health and Care Research; VU
50 University Medical Center; Van der Boechorststraat 7; 1081 BT Amsterdam; The Netherlands;
51 T+31(0)204448392; j.lakerveld@vumc.nl; dedipac@vumc.nl
52
53
54
55
56
57
58
59
60

Abstract

BACKGROUND: The utilisation of available cross-European data for secondary data analyses on physical activity, sedentary behaviours and their underlying determinants may benefit from the wide variation that exists across the continent in terms of these behaviours and their determinants. Such re-use of existing data for further research requires FAIR (Findable; Accessible; Interoperable; Reusable) data management and stewardship. We here describe the inventory and development of a comprehensive European dataset compendium, and the process towards cross European secondary data analyses of pooled data on physical activity, sedentary behaviour and their correlates across the life course.

METHODS: A five-step methodology was followed, covering the 1) identification of relevant datasets across Europe, 2) development of a dataset compendium including details on the design, study population, measures, and level of accessibility of data from each study, 3) definition of key topics and approaches for secondary analyses, 4) process of gaining access to datasets, and 5) pooling and harmonisation of the data and the development of a data harmonisation platform.

RESULTS: A total of 114 unique datasets were found for inclusion within the DEDIPAC compendium. Of these datasets, 14 were eventually obtained and reused to address 10 exemplar research questions. The DEDIPAC data harmonisation platform proved to be useful for pooling, but in general, harmonisation was often restricted to just a few core (crude) outcome variables and some individual-level socio-demographic correlates of these behaviours.

CONCLUSIONS: Obtaining, pooling and harmonising data for secondary data analyses proved to be difficult and sometimes even impossible. Compliance to FAIR data management and stewardship principles currently appears to be limited for research in the field of physical activity and sedentary behaviour. We discuss some of the reasons why this might be the case and present recommendations based on our experience.

Keywords: Physical activity, sedentary behaviour, secondary data analysis, pooling, harmonisation, determinants, Europe, children, adults, older adults

Strengths and Limitations

- We applied the FAIR principles to provide guidance in the discovery and reuse of data for further investigation
- We have identified more than one hundred potentially relevant European datasets through our outlined search approaches
- It was possible to retrieve the metadata from these datasets on the types of variables, age groups under study, study design, measurement instruments used, time frame, etc.
- The reusability of datasets was limited to the small number of datasets to which access was granted
- The variation in assessment methods and operationalization of outcome variables across current European studies hampered data harmonisation process.

Background

Low levels of physical activity and high levels of sedentary behaviour (too much sitting) are recognized as important risk factors for non-communicable disease (NCDs) across Europe.[1–5] These unhealthy behaviours put stress on societies as the burden of NCDs increases among Europe's ageing population.[6] Physical activity and sedentary behaviours are influenced by a wide range of individual-level and contextual factors. Apart from some prominent emerging influencing factors that stand out - such as injuries -, the contribution of most determinants of sedentary behaviour and physical activity patterns may only be relevant for some people, in interaction with other factors, in some places, and under certain circumstances. In other words, such behavioural determinants consist of a combination of individual and contextual factors that may be at the group- or individual level and may mediate and moderate each other. This makes the identification of target factors for policy actions or behavioural interventions challenging. Especially the more distal or 'upstream' behavioural determinants –such as factors in the built or social environments that support or hinder people to be active or sedentary – are often nuanced and hard to identify.[7–9] Sophisticated methods are required for identifying these, using adequately powered data with sufficient variance in physical activity or sedentary behaviours as well as their underlying factors. Gathering such data is a costly endeavour and research funders are increasingly emphasizing the importance of data sharing and secondary analysis of (often publicly funded) datasets as a means to maximize the research potential, increase power and variation of existing resources and provide greater returns on research investments.[10–13]

To address this challenge, one of the aims of the European Determinants of Diet and Physical Activity (DEDIPAC) Knowledge Hub was to explore the interrelations between correlates and determinants of physical activity and sedentary behaviours in and across European populations through secondary data analyses, using state-of-the-art statistical methods.[14] One of the important drivers for establishing the Knowledge Hub was to enable comprehensive and large-scale research using the wide variation in exposures and outcomes that exist across Europe, and to make use of the scattered and sometimes overlapping research that is conducted across its member states.

Utilizing existing cross-European data for secondary data analyses potentially benefits from the wide variety in levels of physical activity and sedentary behaviour and their determinants across the continent.[15–19] Although comparable, objectively measured data is currently lacking, there are indications that adults in northern European countries engage in more sitting time than in countries in the south of Europe,[15] and that some southern European countries generally appear to be among the less physically active countries.[17,19] Analysing combined (i.e., pooled) datasets in which the variables are matched (i.e., harmonised) is a way to increase statistical power, the representativeness of wider populations, and variation in outcomes or correlates of interest. In turn, increasing the power allows for mediation and moderation analyses and subsequent stratified (subgroup) analyses. Such retrospective pooling often requires a data harmonization process in which similar variables across multiple datasets are made compatible.[20] The process of pooling datasets and harmonising the variables can be facilitated by data harmonisation platforms (DHPs). Examples of de-centralised and centralised DHPs that were developed to support the pooling, harmonisation and analyses in epidemiological research include the BioSHaRE-EU [21], the Harvard Dataverse [22] (both de-centralised), and POLARIS [23] (centralised).

1
2
3 Within the numerous organisations and research institutes – inside as well as outside the DEDIPAC
4 consortium – there is a wealth of data on physical activity, sedentary behaviours and their potential
5 individual and contextual correlates, predictors and determinants. Many of these datasets may have
6 strong potential for secondary data analyses if these data become available. It is, however, unclear
7 whether and how these data are distributed across Europe, to what extent these data sets comply to
8 the FAIR principles (i.e. Findable, Accessible, Interoperable, Reusable)[24] for data management and
9 data stewardship, and whether they indeed have the potential to be used in secondary data analyses
10 of behavioural correlates across the life course.
11
12

13
14 In this paper we set out a stepwise process towards cross European data analyses of physical
15 activity, sedentary behaviour and their correlates across the life course. We describe the inventory
16 and development of a comprehensive European dataset compendium as well as the DEDIPAC DHP,
17 and discuss to what extent behavioural physical activity and sedentary behaviour research complies
18 to FAIR principles and can be used for secondary data analyses.
19
20
21
22

23 **Methods**

24
25 The FAIR principles suggest that each data resource, associated metadata and complimentary files
26 should be easy to find ('Findable'); they should provide relevant metadata from these datasets, for
27 instance on the types of variables, age groups under study, study design, measurement instruments
28 used, time frame, etc. ('Accessible'); they should be 'Interoperable' and thus use a consistent data
29 format and taxonomy for knowledge representation and finally, they should be 'Reusable', i.e.,
30 made available. [24]
31
32

33 The DEDIPAC secondary data analysis plan followed a five-step approach including: 1) The
34 identification of relevant datasets, 2) the development of a dataset compendium, 3) the clarification
35 of key topics and approaches for analyses, 4) gaining access to datasets, and 5) pooling of datasets
36 and variable harmonisation. These steps are depicted in Figure 1 and described in further detail
37 below:
38
39

40 1) *Identification of relevant European datasets*

41
42 A relevant European dataset was defined as a dataset collected during an on-going or recently (≤ 10
43 years) completed project focusing on physical activity and/or sedentary behaviour and its potential
44 individual and/or contextual correlates. In addition, further inclusion criteria were formulated as
45 follows: i) Participants had to be 6 years or older, ii) the study had a cross-sectional or longitudinal
46 design, iii) the study was primarily conducted within the European Union, and iv) the dataset
47 consisted of quantitative data which could be either self-reported or objectively measured.
48
49

50 In the context of the DEDIPAC, relevant datasets were deemed 'Findable' if they were identified by:
51 i) A search of the CORDIS project platform, which is the European Commission's primary public
52 repository and portal to disseminate information on all EU-funded research projects and their results
53 in the broadest sense[25], ii) an examination of existing recent reviews of the literature and noting
54 the nature of datasets used, and iii) research network and expert consultation.
55
56
57

58 2) *To develop and complement a compendium of relevant European datasets*

59
60

1
2
3 In an attempt to enhance the findability of the relevant datasets for future searches, a detailed
4 compendium was developed. The compendium is a database in which the following information was
5 detailed: project name, contact person details, website URL (if any), brief description of project,
6 relevant publications, nations involved, sample size, gender, age, physical activity and sedentary
7 behaviour and correlate measurement, indices of inequality/ethnic minorities and level of
8 reusability. This information was gathered from publically available resources. The custodians of the
9 datasets mentioned in the compendium were then approached and asked whether the information
10 was correct and, in some cases, if they could provide additional details. Furthermore, they were
11 asked for their permission to include their project and the project details in the compendium, which
12 would be made accessible to the wider research community and to indicate the level of accessibility
13 for secondary data analysis. The initial willingness of dataset owners to provide the data for re-use
14 was listed in the compendium. This feedback provided some insights into the degree of challenge
15 that would exist in achieving access to targeted datasets.
16
17
18
19

20
21 *3) To define key topics and approaches for analyses*

22 A number of research questions were formulated by the DEDIPAC consortium to assess the potential
23 for secondary data analyses of the identified datasets. We aimed to add to the current state of
24 knowledge as recently systematically summarised[7–9], and informed by the DEDIPAC frameworks
25 on determinants of sedentary behaviour[26] and physical activity. The formulation of research
26 questions was based on three distinct approaches: 1) Clarify linkages of clusters and systems
27 identified in the frameworks, 2) differentiate and nuance correlates of the two behaviours, and 3)
28 begin to fill knowledge gaps in determinant research. An expression of interest (Eoi) statement was
29 requested from DEDIPAC members that were interested in addressing one or more research
30 questions. In addition to a clearly defined research question the Eoi included details of the target
31 population, the project hypothesis, target dataset(s), independent and dependent variables,
32 anticipated data harmonisation approach and the foreseen statistical analysis. To assist in this latter
33 step, a two-day statistical analysis workshop was organised in Amsterdam, the Netherlands,
34 specifically focused on challenges of conducting secondary data analysis, handling pooling and
35 harmonisation issues, and to provide support on advanced statistical techniques (e.g., Bayesian
36 analyses, mediation/moderation analyses and handling missing data).
37
38
39
40
41

42 *4) To co-ordinate access to target datasets*

43 FAIR principles suggest that in order to ensure that existing and future data is accessible and
44 reusable, specific requirements are in place, such as a common and easily shared data format, a
45 common taxonomy, detailed metadata and a data access protocol with a clearly defined data usage
46 licence. Implicit in these requirements are a pathway to or existing ethical approval to share data,
47 safe and secure technologies to transfer data or facilitate remote access and analysis of data, and
48 detailed data dictionaries that clearly define the methodologies used in data collection. In the
49 context of the DEDIPAC project dataset owners of the required/targeted datasets were contacted to
50 assess the potential for accessing datasets and the timeframe and relevant procedure to gain access.
51 After initial and informal consent was obtained, a formal and detailed request for data sharing was
52 provided, including a draft data sharing agreement that covered legal issues, terms of data usage
53 and co-authorships agreements. The precise purpose of the data was specified and a detailed list of
54 variables was requested. After formal agreement the required data was sent to the central data-
55
56
57
58
59
60

1
2
3 managing centre of the data pooling taskforce. The relevant DEDIPAC partners then received these
4 data for analysis after they signed a data sharing agreement for recipients. In the latter agreement,
5 data related issues such as access rights, use, liability, publication and intellectual property rights
6 were formally agreed upon.
7

9 5) *Data pooling and harmonisation of variables*

10
11 The exemplar projects used secondary data analysis on either (i) a single dataset or (ii) pooled and
12 harmonised data from multiple datasets. In the latter case, pooling and harmonising could be
13 conducted manually in a statistical program, or using the DEDIPAC DHP. The DEDIPAC DHP is a
14 Microsoft Access based DHP that is based on a DHP developed for the POLARIS project – a project
15 for individual patient data meta analyses and thus fully dependent on data pooling and
16 harmonization.[23] Within this platform, the original datasets are linked with a reference dataset
17 containing all potentially relevant variables from all individual datasets. If the studies measured and
18 reported the same construct in the same way (e.g., self-reported total physical activity based on the
19 IPAQ-short and reported in minutes per day), these variables were linked to the same variable in the
20 reference dataset. However, if there was a difference in terms of concepts (e.g. total physical activity
21 vs. physical activity in leisure time), measurement (e.g., self-reported vs. objective measurements),
22 or reporting (e.g., minutes per day vs. meeting recommendations) these variables were linked to
23 different variables in the reference dataset. Thus, the reference dataset could contain multiple
24 versions that describe a single variable. In a later step, the multiple versions of a variable could be
25 integrated into a unified measurement scheme. This unified measurement scheme describes a set of
26 variables that have been measured in a similar fashion across a number of datasets, thus having
27 potential for harmonisation. Within DEDIPAC this step involved the identification of specific variables
28 across target datasets which could/would be harmonised and required detailed examination of
29 procedures/methods used by each dataset. Through an iterative review process, by a consensus
30 committee, an agreed approach to the actual harmonisation of the variables was found. The
31 selection process required a balance between uniformity (i.e., exact same question wording and
32 data collection procedures) and the acceptance of a certain level of heterogeneity across datasets
33 (i.e., slightly different wording or procedures). Sharing solutions for data harmonisation across
34 research question groups, regular telephone conferences and shared guidance documents were
35 important aspects of this step.
36
37
38
39
40
41
42
43
44
45

46 **Results**

47 *Identified datasets and the DEDIPAC Compendium*

48
49 A total of 114 unique datasets were identified for inclusion within the DEDIPAC compendium (Figure
50 2). The majority of datasets included in the compendium were found through the CORDIS project
51 platform. Other datasets within the compendium were identified by experts as potentially relevant
52 for inclusion. The compendium is accessible through <https://www.dedipac.eu> and builds upon the
53 present level of findability. Details on the accessibility of these datasets are included in Figure 2. The
54 specific aims for gathering data within each dataset varied, with the majority of included datasets
55 falling under one or more of the following six categories:
56
57
58
59
60

- 1) Development, delivery and evaluation of interventions (n=33);
- 2) National surveys, for example, household surveys (n=17);
- 3) National cohort studies (n=11);
- 4) Assessment and development of policy and research strategies (n=8);
- 5) Studies investigating the link between lifestyle/environmental factors and disease (n=7);
- 6) Studies investigating relations and interactions between health behaviours and health (n=7).

The majority of research projects included in the compendium were funded through the European Commission or other European level (53%), with just under a quarter of the related datasets (23%) supported by national funding. Just under half of datasets included data collected in two or more European countries (47%), with the remainder targeted at individual European countries. The study designs employed within the datasets are summarised in Table 1. In brief, 43 datasets focused on one stage of the life course only; children only (n=13), adolescents only (n=2), adults only (n=19) or older adults only (n=9). Other datasets targeted two or more stages of the life course; with adults and older adults (n=20) being the most frequently combined stages within datasets.

A range of measurement tools were employed across datasets to measure physical activity and/or sedentary behaviour (Table 2). Approximately 41% of datasets included within the compendium used self-report tools to assess physical activity, using questionnaire tools designed specifically for their dataset or other routinely used questionnaires. Within the datasets using self-report tools, 27 (57.4%) used a questionnaire specifically designed for that project, with 10 (21.3%) using the IPAQ. A smaller proportion of studies used self-report proxy measures (1.8%) or a combination of self-report and proxy measures (1.8%) to assess physical activity. Approximately 21% of datasets used self-report tools to assess sedentary behaviour.

Only twenty-four datasets included physical activity (n=24) measured using objective tools, either as standalone or alongside subjective tools, for example, questionnaires. Within the datasets that included objective measures, accelerometers were used to assess physical activity within 16 datasets (66.7%), as well as sensors/ smartphones (20.8%), heart rate monitors (4.2%) and multiple monitors, including sensors and accelerometers (4.2%). Sedentary behaviour was measured using objective tools in 5 of the datasets. Heart rate monitors were used to measure sedentary behaviour in one dataset while accelerometers were used to assess sedentary behaviour in 4 datasets, either on their own (n=2) or in combination with subjective tools (n=2).

The number of datasets reporting factors associated with physical activity and sedentary behaviour are summarised in Tables 3a and 3b with 46% measuring physical activity and 15% measuring sedentary behaviour. These datasets were analysed using the following categories; biological (e.g. sex, gender, health status, ethnicity, etc.), psychological (e.g. intentions, attitudes, self-perception, satisfaction, etc.), behavioural (e.g. lifestyle habits, life events, past experiences, etc.), physical environment (e.g. access, neighbourhood walkability/safety, climate, etc.), socio-cultural (e.g. peer and family support, social expectation, local/national identity, etc.), economic (socio economic status, house ownership, etc.) and policy-related (e.g. promotion initiatives, government policy existence and implementation, etc.). Biological, behavioural and psychological determinants/correlates of physical activity were most frequently reported within included datasets (Table 3a). Behavioural and biological determinants/correlates were also frequently reported sedentary behaviour determinants (Table 3b).

Key topics and approaches for analyses

The EoI process and subsequent refinement resulted in a list of 10 research questions that addressed determinants/correlates across the life course, had a balanced focus on physical activity, sedentary behaviour or both, and required a variety of statistical approaches (e.g., using Bayesian network modelling, Chi-squared automatic interaction detection analysis, etc.). Work teams around each research question were formed. A three-day writing retreat was organised in Ghent, Belgium, to further define approaches, to progress specific questions as well as identify unresolved issues regarding the pooling and harmonisation process.

Reusability of target datasets

The metadata contained within the compendium provide details on the potential level of reusability, and what conditions must be met to facilitate access to the datasets (Figure 2). The compendium provides metadata for all included datasets, regardless of the level of reusability of the data.

Of all 114 identified datasets, 43 were deemed to have potential in answering the 10 exemplar questions posed and subsequently data access requests were made to the dataset owners. Of these, actually, 12 datasets were not accessible during the timeframe of the DEDIPAC project, 10 dataset owners did not respond to our request after a number of attempts, and 20 datasets (47%) agreed that access was possible. At the time of analyses, data had been obtained from 14 dataset owners.

Data Pooling and Harmonisation

The dataset owners that agreed to participate transferred the variables outlined in the data request form. The pooling of these data broadly followed the FAIR principles for interoperability, to facilitate the pooling of data across non-co-operating resources. The data was transferred in a compatible electronic format (most often SPSS, sometimes SAS or STATA). Data was then transferred to SPSS (IBM SPSS Statistics for Windows, Version 22.0. Armonk, NY, USA) and the original data was archived for backup purposes. The datasets were then distributed to those leading the analyses. As a proof of concept, the DEDIPAC DHP was used for one research question. This comprised the pooling of two large datasets and the first simple harmonisation steps encompassed the alignment of gender, education coding, etc.

Discussion

There is an increasing emphasis on the potential utility of existing datasets as a resource to answer research questions. While there are obvious benefits in terms of increased power for statistical analysis, this new and evolving approach has a number of challenges. The DEDIPAC consortium sought, for the first time, to examine the potential for data pooling and secondary data analysis in the measurement and determinants of dietary, physical activity and sedentary behaviours. In so doing, DEDIPAC has developed a valuable compendium of datasets that make it possible to ascertain the scope and scale of European projects in this research domain.

In our approach we applied the FAIR principles to provide guidance in the discovery and reuse of data for further investigation. We have identified more than one hundred potentially relevant European datasets through our outlined search approaches (*'Findability'*). It was possible to retrieve

1
2
3 the metadata from these datasets on the types of variables, age groups under study, study design,
4 measurement instruments used, time frame, etc. (*'Accessibility'*). These metadata were
5 systematically detailed in the DEDIPAC compendium in order to be *'Interoperable'*. These datasets
6 were, and are, being used to address 10 exemplar research questions (*'Reusability'*) via either direct
7 analyses, de-centralised pooling and harmonisation of multiple datasets, or centralised pooling and
8 harmonisation using a DHP. Therefore, we can conclude that, as a proof-of-concept, it is possible to
9 apply the FAIR principles and successfully undertake research projects using existing data.
10
11

12 However, along the way, we encountered a number of significant challenges that were specific to
13 the elements of the FAIR principles, but we also faced issues that applied to the pooling,
14 harmonising and/or analysing the secondary data. These issues should be carefully considered if
15 data pooling and harmonisation are to become a more central aspect of pan-European data analysis.
16
17

18 Firstly, data that does not exist is not possible to find. The inventory of data in the development of
19 the compendium brought to light that the number of existing datasets that may be relevant for
20 pooling, harmonising and secondary data analyses in this field of research was rather limited – or
21 very well hidden. There were especially few current datasets that contained a wide set of variables
22 (outcomes, independent variables, co-variables) to study behaviours and their underlying factors.
23 Therefore, the DEDIPAC project has highlighted a dearth of data on the determinants of physical
24 activity and sedentary behaviour and believe a pan-European cohort study with a focus on
25 behavioural determinants is required to address this deficiency.
26
27
28

29 Secondly, while we have demonstrated that it is feasible to retrieve metadata from datasets, it often
30 required repeated personal (e-mail) contact with data set owners. In addition, and more
31 importantly, the reusability of datasets was limited to the small number of datasets to which access
32 was granted. In the timeframe of the DEDIPAC project it was challenging and time consuming to
33 pursue access to the datasets to address the 10 exemplar research questions and only 12% had
34 provided access at the time of writing. In the future, different approaches may be required to
35 encourage dataset owners to participate in pooled data analysis. The outcome of the DEDIPAC
36 exemplar projects may provide the evidence of the benefits and limitations that would provide
37 dataset owners with a more solid basis for engagement. It is also possible that funding agencies may
38 mandate the sharing of publicly funded data. Therefore, a continued and open discussion of the
39 merits and limitations of data pooling is necessary.
40
41
42

43 Thirdly, and beyond FAIR, the harmonization of core outcome measures under study (in most cases
44 sedentary behaviour and physical activity) was often problematic. Across the included and accessible
45 datasets, these outcomes were measured or operationalized in a variety of ways. The substantial
46 variation in assessment methods and operationalization of outcome variables across current
47 European studies (as illustrated in Table 2) not only hampered the practical harmonisation process,
48 but also presented comparability issues, as estimations of physical activity and sedentary behaviour
49 levels are known to differ based on the assessment method used.[16–19]
50
51
52

53 Fourthly, next to harmonisation issues of core outcomes (physical activity and/or sedentary
54 behaviour, in this context), our focus on determinants of physical activity and sedentary behaviours
55 meant that individual-level and contextual – more upstream – factors were to be taken into account.
56 It was possible to harmonise some of the core outcomes sometimes, and some key socio-
57 demographic variables (e.g. age, gender and educational background) could be harmonised. However,
58
59
60

1
2
3 harmonisation of the behavioural determinants was rarely possible and important co-variables even
4 less so. Thus, pooling often implied that certain variables could not be taken into account if they
5 were not measured across all included datasets, or could not be harmonised. In our opinion, the
6 different assessment methods impeded harmonisation.
7

8
9 Finally, the retrospective pooling and harmonization of variables requires a 'flexible design', since
10 very few established studies have used identical collection methods and procedures.[27] A flexible
11 design means that various categories of a variable such as, for instance, education attainment are
12 reduced to few, or even only two categories (e.g. higher/lower education). Hence, putting data from
13 respondents from different studies together in one dataset indeed increases the number of
14 observations and the absolute power of observations, but the promise that smaller associations can
15 be picked up is likely to be compromised. The added nuance sought for with pooling data may thus
16 be undone or even reversed by the – often rough – re-categorisation of variables during the
17 harmonisation process.
18
19
20
21
22

23 **Conclusions**

24
25 The DEDIPAC project has identified and compiled a large number of pan-European projects related
26 to physical activity and sedentary behaviour. In the design and analysis of 10 exemplar projects,
27 using a variety of approaches, we identified a number of challenges including (i) the limited
28 availability of datasets that contain variables to examine the determinants of physical activity and
29 sedentary behaviour, (ii) the difficulty to establish communication with dataset owners or getting
30 their agreement to share data for analysis and (iii) a low harmonisation potential for the limited
31 number of variables that were available, especially for the potential determinants further upstream.
32 Compliance to FAIR data management and stewardship principles currently appears to be limited for
33 research in the field of physical activity and sedentary behaviour. It is recognized that researchers
34 should be facilitated, funded, requested, or required to share their data and comply to the FAIR
35 principles [13,28]. While recognising the importance of utilising existing data it is equally, if not more
36 important, to highlight the absence of data necessary to investigate the determinants of these
37 behaviours. The lack of suitable data will severely limit the ability of research to inform policy. A
38 bigger, targeted and more standardized data collection is needed in order to maximize the potential
39 of data pooling and harmonisation. Until then, there are only narrow margins for determinant
40 research to build and harvest on previous work.
41
42
43
44
45
46

47 **Declarations**

48 *Competing interests*

49 The authors declare that they have no competing interests.
50
51

52 *Data sharing statement*

53 The data may be obtained from the authors for academic purposes
54
55

56 *Funding*

57
58
59
60

1
2
3 The preparation of this paper was supported by the DEterminants of Diet and Physical Activity
4 (DEDIPAC) knowledge hub. This work is supported by the Joint Programming Initiative 'Healthy Diet
5 for a Healthy Life'. The funding agencies supporting this work are: Belgium: Research Foundation –
6 Flanders; France: Institut National de la Recherche Agronomique (INRA); Germany: Federal Ministry
7 of Education and Research; Italy: Ministry of Education, University and Research (DEDIPAC F.S.
8 02.15.02 COD. B84G14000040008; CDR2.PRIN 2010/11 COD. 2010KL2Y73_003)/ Ministry of
9 Agriculture Food and Forestry Policies; Ireland: The Health Research Board (HRB); The Netherlands:
10 The Netherlands Organisation for Health Research and Development (ZonMw); The United Kingdom:
11 The Medical Research Council (MRC).
12
13

14 15 16 17 **Authors' contributions**

18
19 JB, GC, SC, HvdP, JL, LC and CMD conceived the study. FL and AC assisted in data collection. AL and JK
20 redesigned the DEDIPAC data warehouse. JL drafted the manuscript and AL, FL, MDC, HvdP, DOG,
21 AC, LC, JK, JMO, SC, GC, JB and CMD contributed to iterative revisions.
22
23

24 25 26 **References**

- 27
28 1 WHO. Physical activity Factsheet N°385, Reviewed June 2016.
- 29
30 2 Wilmot EG, Edwardson CL, Achana FA, *et al.* Sedentary time in adults and the association with
31 diabetes, cardiovascular disease and death: systematic review and meta-analysis.
32 *Diabetologia* 2012.
- 33
34 3 Grontved A, Hu F. Television viewing and risk of type 2 diabetes, cardiovascular disease, and
35 all-cause mortality: a meta-analysis. *JAMA* 2011;**305**:2448–55.
- 36
37 4 Biswas A, Oh PI, Faulkner GE, *et al.* Sedentary Time and Its Association With Risk for Disease
38 Incidence, Mortality, and Hospitalization in Adults. *Ann Intern Med* 2015;**162**:123.
39 doi:10.7326/M14-1651
- 40
41 5 Lee IM, Shiroma EJ, Lobelo F, *et al.* Effect of physical inactivity on major non-communicable
42 diseases worldwide: An analysis of burden of disease and life expectancy. *Lancet*
43 2012;**380**:219–29. doi:10.1016/S0140-6736(12)61031-9
- 44
45 6 WHO, Aging NI on, National Institutes of Health. Global Health and Aging. 2011.
- 46
47 7 Stierlin AS, De Lepeleere S, Cardon G, *et al.* A systematic review of determinants of sedentary
48 behaviour in youth: a DEDIPAC-study. *Int J Behav Nutr Phys Act* 2015;**12**:133.
49 doi:10.1186/s12966-015-0291-4
- 50
51 8 O'Donoghue G, Perchoux C, Mensah K, *et al.* A systematic review of correlates of sedentary
52 behaviour in adults aged 18-65 years: a socio-ecological approach. *BMC Public Health*
53 2016;**16**:163. doi:10.1186/s12889-016-2841-3
- 54
55 9 Chastin SFM, Buck C, Freiburger E, *et al.* Systematic literature review of determinants of
56 sedentary behaviour in older adults: a DEDIPAC study. *Int J Behav Nutr Phys Act* 2015;**12**:127.
57 doi:10.1186/s12966-015-0292-3
58
59
60

- 1
2
3 10 Doiron D, Burton P, Marcon Y, *et al.* Data harmonization and federated analysis of
4 population-based studies: the BioSHaRE project. *Emerg Themes Epidemiol* 2013;**10**:12.
5 doi:10.1186/1742-7622-10-12
6
- 7 11 Pisani E, AbouZahr C. Sharing health data: good intentions are not enough. *Bull World Health*
8 *Organ* 2010;**88**:462–6. doi:10.2471/BLT.09.074393
9
- 10 12 Schofield PN, Eppig J, Huala E, *et al.* Sustaining the Data and Bioresource Commons. *Science*
11 *(80-)* 2010;**330**:592–3. doi:10.1126/science.1191506
12
- 13 13 Piwowar HA, Becich MJ, Bilofsky H, *et al.* Towards a data sharing culture: Recommendations
14 for leadership from academic health centers. *PLoS Med.* 2008;**5**:1315–9.
15 doi:10.1371/journal.pmed.0050183
16
- 17 14 Lakerveld J, van der Ploeg HP, Kroeze W, *et al.* Towards the integration and development of a
18 cross-European research network and infrastructure: the DEterminants of Diet and Physical
19 ACTivity (DEDIPAC) Knowledge Hub. *Int J Behav Nutr Phys Act* 2014;**11**:143.
20 doi:10.1186/s12966-014-0143-7
21
- 22 15 Loyen A, Van Der Ploeg HP, Bauman A, *et al.* European sitting championship: Prevalence and
23 correlates of self-reported sitting time in the 28 European Union Member States. *PLoS One*
24 2016;**11**.
25
- 26 16 Loyen A, Verloigne M, Van Hecke L, *et al.* Variation in population levels of sedentary time in
27 European adults according to cross-European studies: a systematic literature review within
28 DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**:1–11. doi:10.1186/s12966-016-0397-3
29
- 30 17 Van Hecke L, Loyen A, Verloigne M, *et al.* Variation in population levels of physical activity in
31 European children and adolescents according to cross-European studies: a systematic
32 literature review within DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**:1–22.
33 doi:10.1186/s12966-016-0396-4
34
- 35 18 Verloigne M, Loyen A, Van Hecke L, *et al.* Variation in population levels of sedentary time in
36 European children and adolescents according to cross-European studies: a systematic
37 literature review within DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**:1–30.
38 doi:10.1186/s12966-016-0395-5
39
- 40 19 Loyen A, Van Hecke L, Verloigne M, *et al.* Variation in population levels of physical activity in
41 European adults according to cross-European studies: a systematic literature review within
42 DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**.
43
- 44 20 Granda P, Blasczyk E. Data harmonization, guidelines for best practice in cross-cultural
45 surveys. *MI Surv Res Centre, Inst Soc Res Univ Michigan* 2010.
46
- 47 21 Gaye A, Marcon Y, Isaeva J, *et al.* DataSHIELD: Taking the analysis to the data, not the data to
48 the analysis. *Int J Epidemiol* 2014;**43**:1929–44. doi:10.1093/ije/dyu188
49
- 50 22 Crosas M. The dataverse network®: An open-source application for sharing, discovering and
51 preserving data. *D-Lib Mag* 2011;**17**. doi:10.1045/january2011-crosas
52
- 53 23 Buffart LM, Kalter J, Chinapaw MJM, *et al.* Predicting Optimal cAncer Rehabilitation and
54 Supportive care (POLARIS): rationale and design for meta-analyses of individual patient data
55 of randomized controlled trials that evaluate the effect of physical activity and psychosocial
56 interventions on health. *Syst Rev* 2013;**2**:1–9. doi:10.1186/2046-4053-2-75
57
58
59
60

- 1
2
3 24 Wilkinson MD, Dumontier M, Aalbersberg IJ, *et al.* The FAIR Guiding Principles for scientific
4 data management and stewardship. *Sci Data* 2016;**3**:160018. doi:10.1038/sdata.2016.18
5
6 25 European Commission. CORDIS (Community Research and Development Information
7 Services.
8
9 26 Chastin SFM, De Craemer M, Lien N, *et al.* The SOS-framework (Systems of Sedentary
10 behaviours): an international consensus transdisciplinary framework and research priorities
11 for the study of determinants of sedentary behaviour and policy across the life course: a
12 DEDIPAC study. *Int J BehavNutrPhysAct* 2016;**13**.
13
14 27 Doiron D, Parminder R, Ferretti V, *et al.* Facilitating collaborative research: Implementing a
15 platform supporting data harmonization and pooling. *Nor Epidemiol* 2012;**21**:221–4.
16
17 28 Genbank A. Let data speak to data. *Nature* 2005;**438**:531. doi:10.1038/438531a
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Tables and Figures

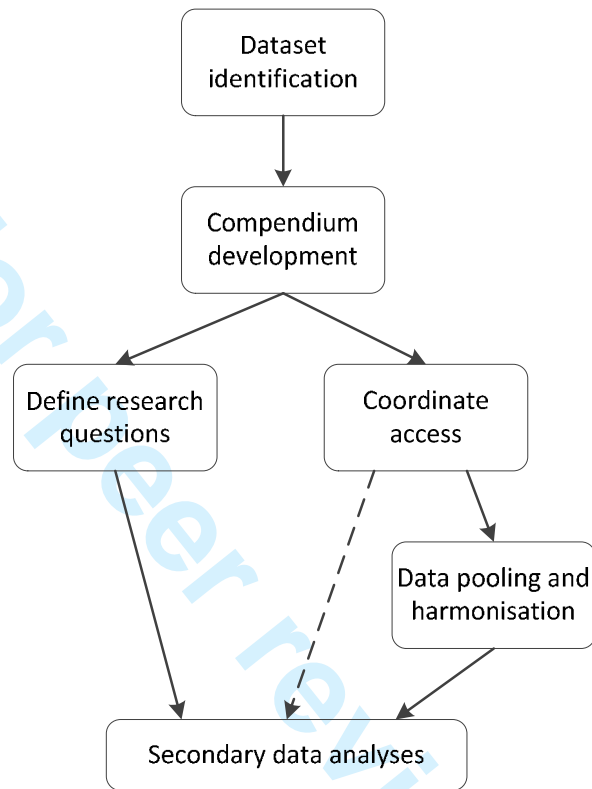
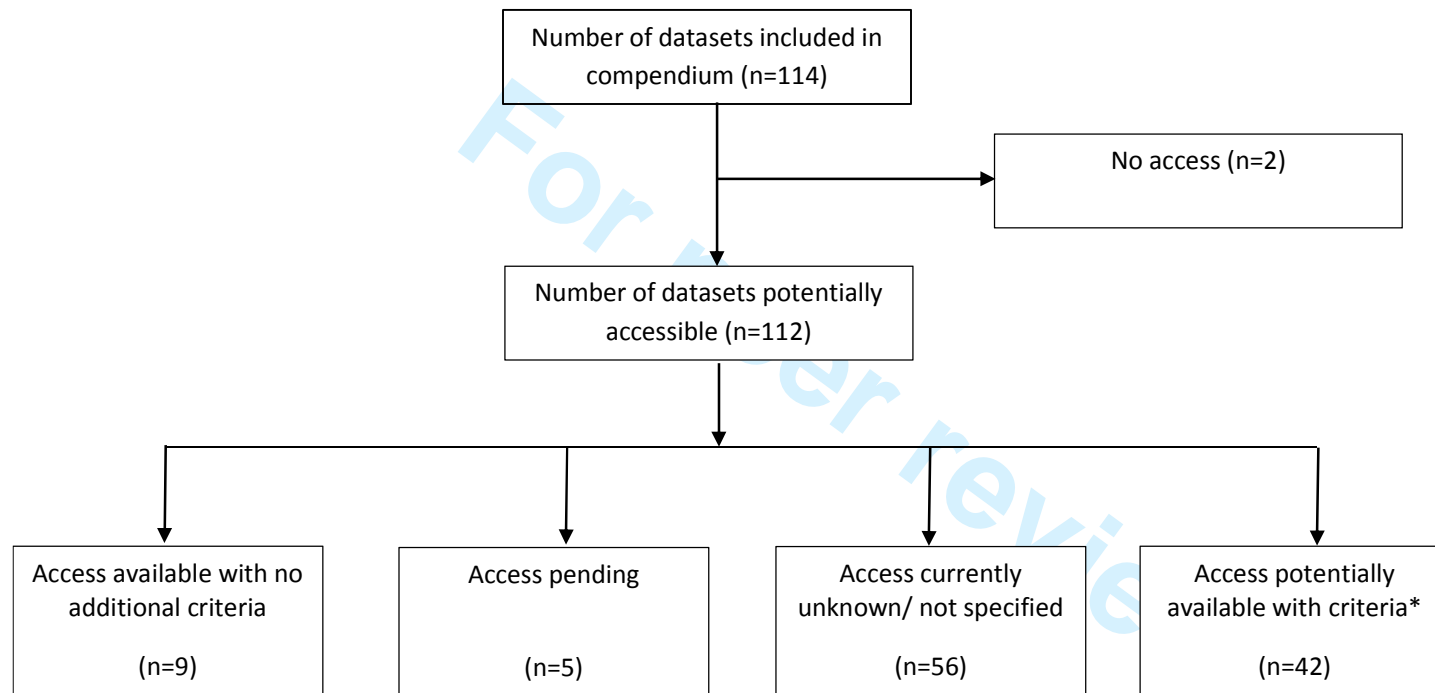
Figure 1: Schematic outline of the stepwise approach towards secondary data analyses

Figure 2: Flow diagram of accessibility to datasets



* Criteria applies regarding access to datasets, for example, Access with Permission, Access after Board Review and Written Material Transfer Agreement, Access to Raw Data at Local Site and/or Provision of Summary Tables for Meta-Analysis (note: for some datasets multiple criteria apply)

Table 1: Nature of datasets included within the DEDIPAC Compendium

	n (%)
Cross-sectional	41 (36)
Longitudinal (including cohort)	17 (15)
Intervention	18 (16)
Cross-sectional and longitudinal	13 (11)
Cross-sectional and intervention	7 (6)
Longitudinal & intervention	2 (2)
Cross-sectional, longitudinal and intervention	2 (2)
Not specified	14 (12)
Total	114

Table 2: Methods used to assess physical activity and/ or sedentary behaviour within datasets

	Physical Activity ^a		Sedentary Behaviour ^a	
	Datasets		Datasets	
	n	%	n	%
Self-report	47	41.2	24	21.1
<i>Tool specifically designed for study</i>	27	57.4	17	70.8
<i>IPAQ</i>	10	21.3	3	12.5
<i>PAQ-S</i>	1	2.1		
<i>GPAQ</i>	2	4.3		
<i>SQUASH</i>	3	6.4		
<i>MAQ</i>	1	2.1		
<i>LAPAQ</i>	1	2.1		
<i>PAQ-A & IPAQ</i>	1	2.1		
<i>IPAQ, RPAQ & EPIC-PAQ</i>	1	2.1	1	4.2
<i>Marshall Questionnaire</i>			1	4.2
<i>AQuAA</i>			2	8.3
Self-report (proxy)	3	2.6	3	2.6
<i>Tool specifically designed for study</i>	3	100	3	100
Self-report and proxy	1	0.9	1	0.9
<i>Tool specifically designed for study</i>	1	100		
<i>IPAQ</i>			1	100
Objective	24	21.1	5	4.4
<i>Sensors/smartphones</i>	5	20.8		
<i>Accelerometer</i>	16	66.7	4	80
<i>Heart rate monitor</i>	1	4.2	1	20
<i>Multiple monitors *</i>	1	4.2		

* SenseWear Armband, Kenz, Actigraph GT3X, Dynaport MiniMod^a Sub category percentage calculation is a percentage of the total number of databases for each measurement method. *Note: some datasets may have used more than one method to assess physical activity and/ or sedentary behaviour. IPAQ, International Physical Activity Questionnaire; PAQ-S, Physical Activity Questionnaire for Schoolchildren; GPAQ, Global Physical Activity Questionnaire; SQUASH, Short Questionnaire to assess health-enhancing physical activity; MAQ, Modifiable Activity Questionnaire; LAPAQ, LASA Physical Activity Questionnaire; PAQ-A, Physical Activity Questionnaire for Adolescents; RPAQ, Recent Physical Activity Questionnaire; EPIC-PAQ, EPIC Physical Activity Questionnaire; AQuAA, Activity Questionnaire for Adults and Adolescents*

Table 3a: Breakdown of datasets reporting on physical activity correlates or determinants by stage of the life course (where applicable)

	Datasets including physical activity determinants (overall)	Biological	Psychological	Behavioural	Physical Environmental	Socio- cultural	Economic	Policy
Children only	8	6	2	6	2	6	1	3
Adolescents only	1	1	1	1	1	1	1	0
Adults only	11	7	7	3	3	1	1	0
Older Adults only	3	2	0	1	0	1	0	0
Children and Adolescents	5	2	3	3	1	0	0	0
All stages of life course	2	0	0	1	2	0	0	0
Children, Adults, Older Adults	3	2	2	2	1	0	0	0
Adolescents, Adults, Older Adults	3	1	0	1	2	2	1	0
Adults and Older Adults	11	7	4	5	2	2	1	0
Children and Adults	1	1	0	1	0	0	0	0
Children, Adolescents, Adults	3	2	3	1	2	2	1	0
Adolescents and Adults	2	0	1	1	0	0	0	0
Total	53	31	23	26	16	15	6	3

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Table 3b: Breakdown of datasets reporting on sedentary correlates or determinants by stage of the life course (where applicable)

	Datasets including sedentary behaviour determinants (overall)	Biological	Psychological	Behavioural	Physical Environmental	Socio- cultural	Economic	Policy
Children only	5	1	0	4	0	1	1	0
Adults only	1	1	1	0	1	1	1	0
Children and Adolescents	3	0	0	3	0	0	0	0
All stages of life course	2	0	0	2	0	0	0	0
Children, Adults, Older Adults	2	2	1	0	0	0	1	0
Adolescents, Adults, Older Adults	1	1	0	0	1	1	1	0
Adults and Older Adults	2	1	1	0	1	0	0	0
Children, Adolescents, Adults	1	1	1	0	0	0	1	0
Total	17	7	4	9	3	3	5	0

For peer review only

BMJ Open

Identifying and sharing data for secondary data analysis of physical activity, sedentary behaviour and their determinants across the life course in Europe: general principles and an example from DEDIPAC

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2017-017489.R1
Article Type:	Research
Date Submitted by the Author:	17-Aug-2017
Complete List of Authors:	Lakerveld, Jeroen; VU University Medical Center, Loyen, A; VU University Medical Center, Ling, Fiona; University of Limerick, Physical Education and Sport Sciences Department DeCraemer, Marieke; Ghent University, movement and sport science van der Ploeg, Hidde; VU University Medical Center Amsterdam, Department of Public and Occupational Health; Prevention Research Collaboration, Sydney School of Public Health, University of Sydney O'Gorman, Donal; Dublin City University, Centre for Preventive Medicine Carlin, Angela; University of Limerick, Centre for physical activity and health research Caprinica, Laura; University of Rome Foro Italico Kalter, Joeri; VU University Medical Center Oppert, Jean-Michel; University Pierre et Marie Curie, Institute of cardiometabolism and nutrition (ICAN), Pitie-Salpêtrière hospital (AP-HP) Chastin, Sebastian; Glasgow Caledonian University, Institute of Applied Health Research, School of Health and Life Sciences; Ghent University, Department of Movement and Sports Sciences Cardon, Greet; Ghent University, Movement and Sports Sciences Brug, Johannes; VU University Medical Center , EMGO+, Institute for Health and Care Research MacDonncha, Ciaran; University of Limerick, Centre for Physical Activity and Health Research
Primary Subject Heading:	Epidemiology
Secondary Subject Heading:	Research methods, Public health
Keywords:	EPIDEMIOLGY, PREVENTIVE MEDICINE, PUBLIC HEALTH

SCHOLARONE™
Manuscripts

1
2
3 Identifying and sharing data for secondary data analysis of physical activity,
4 sedentary behaviour and their determinants across the life course in Europe:
5 general principles and an example from DEDIPAC
6
7

8 Jeroen Lakerveld^{1§}, Anne Loyen¹, Fiona Chun Man Ling^{2,3,4}, Marieke De Craemer⁵, Hidde P. van der
9 Ploeg⁶, Donal J. O’Gorman⁷, Angela Carlin², Laura Capranica⁸, Joeri Kalter¹, Jean-Michel Oppert⁹,
10 Sebastien Chastin¹⁰, Greet Cardon⁵, Johannes Brug^{1,11}, Ciaran MacDonncha²
11
12

13
14
15 ¹ Department of Epidemiology and Biostatistics, EMGO Institute for Health and Care Research, VU
16 University medical center, De Boelelaan 1089a, 1081, HV, Amsterdam, The Netherlands.

17
18 ² Centre for Physical Activity and Health Research, Department of Physical Education and Sport
19 Sciences, University of Limerick, Ireland.
20

21
22 ³ Institute of Sport, Exercise and Active Living, Victoria University, Melbourne, Australia.

23
24 ⁴ Department of Psychology, Bournemouth University, Bournemouth, UK

25
26 ⁵ Department of Movement and Sports Sciences, Ghent University, Ghent, Belgium.

27
28 ⁶ Department of Public and Occupational Health, EMGO Institute for Health and Care Research, VU
29 University medical center, Amsterdam, The Netherlands.
30

31
32 ⁷ Centre for Preventive Medicine, School of Health & Human Performance, Dublin City University,
33 Dublin 9, Republic of Ireland.

34
35 ⁸ Department of Movement, Human and Health Sciences, University of Rome Foro Italico, Rome,
36 Italy.

37
38 ⁹ Institute of Cardiometabolism and Nutrition (ICAN), University Pierre et Marie Curie-Paris6;
39 Department of Nutrition, Pitie-Salpetriere hospital (AP-HP), Paris, France.

40
41 ¹⁰ Institute of Applied Health Research, School of Health and Life Science, Glasgow Caledonian
42 University, Glasgow, UK.
43

44
45 ¹¹ Amsterdam School for Communication Research, University of Amsterdam, Amsterdam, the
46 Netherlands.
47

48
49
50
51 [§] **Corresponding author:** Jeroen Lakerveld; EMGO Institute for Health and Care Research; VU
52 University Medical Center; Van der Boechorststraat 7; 1081 BT Amsterdam; The Netherlands;
53 T+31(0)204448392; j.lakerveld@vumc.nl; dedipac@vumc.nl
54
55
56
57
58
59
60

Abstract

BACKGROUND: The utilisation of available cross-European data for secondary data analyses on physical activity, sedentary behaviours and their underlying determinants may benefit from the wide variation that exists across Europe in terms of these behaviours and their determinants. Such re-use of existing data for further research requires FAIR (Findable; Accessible; Interoperable; Reusable) data management and stewardship. We here describe the inventory and development of a comprehensive European dataset compendium, and the process towards cross-European secondary data analyses of pooled data on physical activity, sedentary behaviour and their correlates across the life course.

METHODS: A five-step methodology was followed by the European Determinants of Diet and Physical Activity (DEDIPAC) Knowledge Hub, covering the 1) identification of relevant datasets across Europe, 2) development of a compendium including details on the design, study population, measures, and level of accessibility of data from each study, 3) definition of key topics and approaches for secondary analyses, 4) process of gaining access to datasets, and 5) pooling and harmonisation of the data and the development of a data harmonisation platform.

RESULTS: A total of 114 unique datasets were found for inclusion within the DEDIPAC compendium. Of these datasets, 14 were eventually obtained and reused to address 10 exemplar research questions. The DEDIPAC data harmonisation platform proved to be useful for pooling, but in general, harmonisation was often restricted to just a few core (crude) outcome variables and some individual-level socio-demographic correlates of these behaviours.

CONCLUSIONS: Obtaining, pooling and harmonising data for secondary data analyses proved to be difficult and sometimes even impossible. Compliance to FAIR data management and stewardship principles currently appears to be limited for research in the field of physical activity and sedentary behaviour. We discuss some of the reasons why this might be the case and present recommendations based on our experience.

Keywords: Physical activity, sedentary behaviour, secondary data analysis, pooling, harmonisation, determinants, Europe, children, adults, older adults

Strengths and Limitations

- We applied the FAIR principles to provide guidance in the discovery and reuse of data for further investigation
- We have identified more than one hundred potentially relevant European datasets through our outlined search approaches
- It was possible to retrieve the metadata from these datasets on the types of variables, age groups under study, study design, measurement instruments used, time frame, etc.
- Limited potential for reuse has been noted and this highlights the immediate need to manage future data collection within Europe using the FAIR principles

- More consistent data collection methodologies among the scientific community should be promoted as the variation in assessment methods and operationalization of outcome variables across current European studies hampered data harmonisation.

Background

Low levels of physical activity and high levels of sedentary behaviour (too much sitting) are recognized as important risk factors for non-communicable disease (NCDs) across Europe.[1–5] These unhealthy behaviours put stress on societies as the burden of NCDs increases among Europe’s ageing population.[6] Physical activity and sedentary behaviours are influenced by a wide range of individual-level and contextual factors.[7–13] Apart from some prominent emerging influencing factors that stand out - such as injuries -, the contribution of most determinants of sedentary behaviour and physical activity patterns may only be relevant for some people, in interaction with other factors, in some places, and under certain circumstances. In other words, such behavioural determinants consist of a combination of individual and contextual factors that may be at the group- or individual level and may mediate and moderate each other. This makes the identification of target factors for policy actions or behavioural interventions challenging. Especially the more distal or ‘upstream’ behavioural determinants –such as factors in the built or social environments that support or hinder people to be active or sedentary – are often nuanced and hard to identify.[11,13,12,14] Sophisticated methods are required for identifying these, using adequately powered data with sufficient variance in physical activity or sedentary behaviours as well as their underlying factors. Gathering such data is a costly endeavour and research funders are increasingly emphasizing the importance of data sharing and secondary analysis of (often publicly funded) datasets as a means to maximize the research potential, increase power and variation of existing resources and provide greater returns on research investments.[15–18]

To address this challenge, one of the aims of the European Determinants of Diet and Physical Activity (DEDIPAC) Knowledge Hub was to explore the interrelations between correlates and determinants of physical activity and sedentary behaviours in and across European populations through secondary data analyses, using state-of-the-art statistical methods.[19] One of the important drivers for establishing the Knowledge Hub was to enable comprehensive and large-scale research using the wide variation in exposures and outcomes that exist across Europe, and to make use of the scattered and sometimes overlapping research that is conducted across its member states.

Utilizing existing cross-European data for secondary data analyses potentially benefits from the wide variety in levels of physical activity and sedentary behaviour and their determinants across the continent.[20–24] Although comparable, objectively measured data are currently lacking, there are indications that adults in northern European countries engage in more sitting time than in countries in the south of Europe,[20] and that some southern European countries generally appear to be among the less physically active countries.[22,24] Analysing combined (i.e., pooled) datasets in which the variables are matched (i.e., harmonised) is a way to increase statistical power, the representativeness of wider populations, and variation in outcomes or correlates of interest. In turn, increasing the power allows for mediation and moderation analyses and subsequent stratified (subgroup) analyses. Also the re-analysis of previously collected data using contemporary statistical techniques is possible. Such retrospective pooling often requires a data harmonization process in

1
2
3 which similar variables across multiple datasets are made compatible.[25] The process of pooling
4 datasets and harmonising the variables can be facilitated by data harmonisation platforms (DHPs).
5 Examples of de-centralised and centralised DHPs that were developed to support the pooling,
6 harmonisation and analyses in epidemiological research include the BioSHaRE-EU [26], the Harvard
7 Dataverse [27] (both de-centralised), and ICAD [28] and POLARIS [29] (both centralised).
8
9

10 Within the numerous organisations and research institutes – inside as well as outside the DEDIPAC
11 consortium – there is a wealth of data on physical activity, sedentary behaviours and their potential
12 individual and contextual correlates, predictors and determinants. Many of these datasets may have
13 strong potential for secondary data analyses if these data become available. It is, however, unclear
14 whether and how these data are distributed across Europe, to what extent these datasets comply to
15 the FAIR principles (i.e. Findable, Accessible, Interoperable, Reusable)[30] for data management and
16 data stewardship, and whether they indeed have the potential to be used in secondary data analyses
17 of behavioural correlates across the life course.
18
19

20
21 In this paper we set out a stepwise process towards cross-European data analyses of physical
22 activity, sedentary behaviour and their correlates across the life course. We describe the inventory
23 and development of a comprehensive European dataset compendium as well as the DEDIPAC DHP,
24 and discuss to what extent behavioural physical activity and sedentary behaviour research complies
25 to FAIR principles and can be used for secondary data analyses.
26
27
28
29

30 **Methods**

31
32 The FAIR principles suggest that each data resource, associated metadata and complementary files
33 should be registered or indexed in a searchable resource, so that they can be located ('Findable');
34 they should provide relevant metadata from these datasets, for instance on the types of variables,
35 age groups under study, study design, measurement instruments used, time frame, etc.
36 ('Accessible'); they should be 'Interoperable' and thus use a consistent data format and taxonomy
37 for knowledge representation and finally, they should be 'Reusable', i.e., made available. [30]
38
39

40 The DEDIPAC secondary data analysis plan followed a five-step approach including: 1) The
41 identification of relevant datasets, 2) the development of a dataset compendium, 3) the clarification
42 of key topics and approaches for analyses, 4) gaining access to datasets, and 5) pooling of datasets
43 and variable harmonisation. These steps are depicted in Figure 1 and described in further detail
44 below:
45
46

47 1) *Identification of relevant European datasets*

48
49 A relevant European dataset was defined as a dataset collected during an on-going or recently (≤ 10
50 years) completed observational or intervention study focusing on physical activity and/or sedentary
51 behaviour and its potential individual and/or contextual correlates. In addition, further inclusion
52 criteria were formulated as follows: i) Participants had to be aged 6 years or older, ii) the
53 observational or intervention study had a cross-sectional or longitudinal design, iii) the study was
54 primarily conducted within the European Union, and iv) the dataset consisted of quantitative data
55 which could be either self-reported or objectively measured.
56
57
58
59
60

1
2
3 In the context of the DEDIPAC, relevant datasets were deemed 'Findable' if they were identified by:
4 i) A search of the CORDIS project platform, which is the European Commission's primary public
5 repository and portal to disseminate information on all EU-funded research projects and their results
6 in the broadest sense[31], ii) an examination of existing recent reviews of the literature and noting
7 the nature of datasets used, and iii) activities of the DEDIPAC Knowledge Hub and expert
8 consultation.
9

10
11
12 2) *To develop and complement a compendium of relevant European datasets*

13 In an attempt to enhance the findability of the relevant datasets for future searches, a detailed
14 compendium was developed. The compendium is a database in which the following information was
15 detailed: project name, contact person details, website URL (if any), brief description of project,
16 relevant publications, nations involved, sample size, gender, age, physical activity and sedentary
17 behaviour and correlate measurement, indices of inequality/ethnic minorities and level of
18 reusability. This information was gathered from publically available resources. The custodians of the
19 datasets mentioned in the compendium were then approached and asked whether the information
20 was correct and, in some cases, if they could provide additional details. Furthermore, they were
21 asked for their permission to include their project and the project details in the compendium, which
22 would be made accessible to the wider research community and to indicate the level of accessibility
23 for secondary data analysis. The initial willingness of dataset owners to provide the data for re-use
24 was listed in the compendium. This feedback provided some insights into the degree of challenge
25 that would exist in achieving access to targeted datasets.
26
27
28
29
30

31 3) *To define key topics and approaches for analyses*

32 A number of research questions were formulated by the DEDIPAC consortium to assess the potential
33 for secondary data analyses of the identified datasets. We aimed to add to the current state of
34 knowledge as recently systematically summarised[11,13,12], and informed by the DEDIPAC
35 frameworks on determinants of sedentary behaviour[32] and physical activity. The formulation of
36 research questions was based on three distinct approaches: 1) Clarify linkages of clusters and
37 systems identified in the frameworks, 2) differentiate and nuance correlates of the two behaviours,
38 and 3) begin to fill knowledge gaps in determinant research. An expression of interest (Eoi)
39 statement was requested from DEDIPAC members that were interested in addressing one or more
40 research questions. In addition to a clearly defined research question the Eoi included details of the
41 target population, the project hypothesis, target dataset(s), independent and dependent variables,
42 anticipated data harmonisation approach and the foreseen statistical analysis. To assist in this latter
43 step, a two-day statistical analysis workshop was organised in Amsterdam, the Netherlands,
44 specifically focused on challenges of conducting secondary data analysis, handling pooling and
45 harmonisation issues, and to provide support on advanced statistical techniques (e.g., Bayesian
46 analyses, mediation/moderation analyses and handling missing data).
47
48
49
50
51

52 4) *To co-ordinate access to target datasets*

53 FAIR principles suggest that in order to ensure that existing and future data are accessible and
54 reusable, specific requirements are in place, such as a common and easily shared data format, a
55 common taxonomy, detailed metadata and a data access protocol with a clearly defined data usage
56 licence. Implicit in these requirements is a pathway to or existing ethical approval to share data, safe
57
58
59
60

1
2
3 and secure technologies to transfer data or facilitate remote access and analysis of data, and
4 detailed data dictionaries that clearly define the methodologies used in data collection. In the
5 context of the DEDIPAC project dataset owners of the required/targeted datasets were contacted to
6 assess the potential for accessing datasets and the timeframe and relevant procedure to gain access.
7 After initial and informal consent was obtained, a formal and detailed request for data sharing was
8 provided, including a draft data sharing agreement that covered legal issues, terms of data usage
9 and co-authorships agreements. The precise purpose of the data was specified and a detailed list of
10 variables was requested. After formal agreement the required data were sent to the central data-
11 managing centre of the data pooling taskforce. The relevant DEDIPAC partners then received these
12 data for analysis after they signed a data sharing agreement for recipients. In the latter agreement,
13 data related issues such as access rights, use, liability, publication and intellectual property rights
14 were formally agreed upon.
15
16
17
18

19 5) *Data pooling and harmonisation of variables*

20
21 The exemplar projects used secondary data analysis on either (i) a single dataset or (ii) pooled and
22 harmonised data from multiple datasets. In the latter case, pooling and harmonising could be
23 conducted manually in a statistical program, or using the DEDIPAC DHP. The DEDIPAC DHP is a
24 Microsoft Access based DHP that is based on a DHP developed for the POLARIS project – a project
25 for individual patient data-meta analyses and thus fully dependent on data pooling and
26 harmonization.[29] Within this platform, the original datasets are linked with a reference dataset
27 containing all potentially relevant variables from all individual datasets. If the studies measured and
28 reported the same construct in the same way (e.g., self-reported total physical activity based on the
29 IPAQ-short and reported in minutes per day), these variables were linked to the same variable in the
30 reference dataset. However, if there was a difference in terms of concepts (e.g. total physical activity
31 vs. physical activity in leisure time), measurement (e.g., self-reported vs. objective measurements),
32 or reporting (e.g., minutes per day vs. meeting recommendations) these variables were linked to
33 different variables in the reference dataset. Thus, the reference dataset could contain multiple
34 variables for a single construct. In a later step, the multiple versions of a variable could be integrated
35 into a unified measurement scheme. This unified measurement scheme describes a set of variables
36 that have been measured in a similar fashion across a number of datasets, thus having potential for
37 harmonisation. Within DEDIPAC this step involved the identification of specific variables across
38 target datasets which could/would be harmonised and required detailed examination of
39 procedures/methods used by each dataset. Through an iterative review process, by a consensus
40 committee, an agreed approach to the actual harmonisation of the variables was found. The
41 selection process required a balance between uniformity (e.g., exact same question wording and
42 data collection procedures, or using the same data processing standards and cut-points) and the
43 acceptance of a certain level of heterogeneity across datasets (i.e., slightly different wording or
44 procedures). As an aspect of expert consensus and where possible an examination of the frequency
45 and descriptive statistics of similar variables informed the potential for harmonisation. Sharing
46 solutions for data harmonisation across research question groups, regular telephone conferences
47 and shared guidance documents were important aspects of this step.
48
49
50
51
52
53
54
55
56
57
58
59
60

Results

Identified datasets and the DEDIPAC Compendium

A total of 114 unique datasets were identified for inclusion within the DEDIPAC compendium (Figure 2). The majority of datasets included in the compendium were found through the CORDIS project platform. Other datasets within the compendium were identified by experts as potentially relevant for inclusion. The compendium is accessible through <https://www.dedipac.eu> and builds upon the present level of findability. Details on the accessibility of these datasets are included in Figure 2. The specific aims for gathering data within each dataset varied, with the majority of included datasets falling under one or more of the following six categories:

- 1) Development, delivery and evaluation of interventions (n=33);
- 2) National surveys, for example, household surveys (n=17);
- 3) National cohort studies (n=11);
- 4) Assessment and development of policy and research strategies (n=8);
- 5) Studies investigating the link between lifestyle/environmental factors and disease (n=7);
- 6) Studies investigating relations and interactions between health behaviours and health (n=7).

The majority of research projects included in the compendium were funded through the European Commission or other European level (53%), with just under a quarter of the related datasets (23%) supported by national funding. Just under half of datasets included data collected in two or more European countries (47%), with the remainder targeted at individual European countries. The study designs employed within the datasets are summarised in Table 1. In brief, 43 datasets focused on one stage of the life course only; children only (n=13), adolescents only (n=2), adults only (n=19) or older adults only (n=9). Other datasets targeted two or more stages of the life course; with adults and older adults (n=20) being the most frequently combined stages within datasets.

Table 1: Nature of datasets included within the DEDIPAC Compendium

	n (%)
Cross-sectional	41 (36)
Longitudinal (including cohort)	17 (15)
Intervention	18 (16)
Cross-sectional and longitudinal	13 (11)
Cross-sectional and intervention	7 (6)
Longitudinal & intervention	2 (2)
Cross-sectional, longitudinal and intervention	2 (2)
Not specified	14 (12)
Total	114

A range of measurement tools were employed across datasets to measure physical activity and/or sedentary behaviour (Table 2). Approximately 41% of datasets included within the compendium used self-report tools to assess physical activity, using questionnaire tools designed specifically for their dataset or other routinely used questionnaires. Within the datasets using self-report tools, 27

(57.4%) used a questionnaire specifically designed for that project, with 10 (21.3%) using the IPAQ. A smaller proportion of studies used self-report proxy measures (1.8%) or a combination of self-report and proxy measures (1.8%) to assess physical activity. Approximately 21% of datasets used self-report tools to assess sedentary behaviour.

Only twenty-four datasets included physical activity (n=24) measured using objective tools, either as standalone or alongside subjective tools, for example, questionnaires. Within the datasets that included objective measures, accelerometers were used to assess physical activity within 16 datasets (66.7%), as well as sensors/ smartphones (20.8%), heart rate monitors (4.2%) and multiple monitors, including sensors and accelerometers (4.2%). Sedentary behaviour was measured using objective tools in 5 of the datasets. Heart rate monitors were used to measure sedentary behaviour in one dataset while accelerometers were used to assess sedentary behaviour in 4 datasets, either on their own (n=2) or in combination with subjective tools (n=2).

Table 2: Methods used to assess physical activity and/ or sedentary behaviour within datasets

	Physical Activity ^a		Sedentary Behaviour ^a	
	Datasets		Datasets	
	n	%	n	%
Self-report	47	41.2	24	21.1
<i>Tool specifically designed for study</i>	27	57.4	17	70.8
<i>IPAQ</i>	10	21.3	3	12.5
<i>PAQ-S</i>	1	2.1		
<i>GPAQ</i>	2	4.3		
<i>SQUASH</i>	3	6.4		
<i>MAQ</i>	1	2.1		
<i>LAPAQ</i>	1	2.1		
<i>PAQ-A & IPAQ</i>	1	2.1		
<i>IPAQ, RPAQ & EPIC-PAQ</i>	1	2.1	1	4.2
<i>Marshall Questionnaire</i>			1	4.2
<i>AQuAA</i>			2	8.3
Self-report (proxy)	3	2.6	3	2.6
<i>Tool specifically designed for study</i>	3	100	3	100
Self-report and proxy	1	0.9	1	0.9
<i>Tool specifically designed for study</i>	1	100		
<i>IPAQ</i>			1	100
Objective	24	21.1	5	4.4
<i>Sensors/smartphones</i>	5	20.8		
<i>Accelerometer</i>	16	66.7	4	80
<i>Heart rate monitor</i>	1	4.2	1	20
<i>Multiple monitors *</i>	1	4.2		

* SenseWear Armband, Kenz, Actigraph GT3X, Dynaport MiniMod^a Sub category percentage calculation is a percentage of the total number of databases for each measurement method. *Note: some datasets may have used more than one method to assess physical activity and/ or sedentary*

1
2
3 *behaviour. IPAQ, International Physical Activity Questionnaire; PAQ-S, Physical Activity Questionnaire*
4 *for Schoolchildren; GPAQ, Global Physical Activity Questionnaire; SQUASH, Short Questionnaire to*
5 *assess health-enhancing physical activity; MAQ, Modifiable Activity Questionnaire; LAPAQ, LASA*
6 *Physical Activity Questionnaire; PAQ-A, Physical Activity Questionnaire for Adolescents; RPAQ, Recent*
7 *Physical Activity Questionnaire; EPIC-PAQ, EPIC Physical Activity Questionnaire; AQuAA, Activity*
8 *Questionnaire for Adults and Adolescents*
9

10
11
12
13 The number of datasets reporting factors associated with physical activity and sedentary behaviour
14 are summarised in Tables 3a and 3b with 46% measuring physical activity and 15% measuring
15 sedentary behaviour. These datasets were analysed using the following categories; biological (e.g.
16 sex, gender, health status, ethnicity, etc.), psychological (e.g. intentions, attitudes, self-perception,
17 satisfaction, etc.), behavioural (e.g. lifestyle habits, life events, past experiences, etc.), physical
18 environment (e.g. access, neighbourhood walkability/safety, climate, etc.) , socio-cultural (e.g. peer
19 and family support, social expectation, local/national identity, etc.), economic (socio economic
20 status, house ownership, etc.) and policy-related (e.g. promotion initiatives, government policy
21 existence and implementation, etc.). Biological, behavioural and psychological
22 determinants/correlates of physical activity were most frequently reported within included datasets
23 (Table 3a). Behavioural and biological determinants/correlates were also frequently reported
24 sedentary behaviour determinants (Table 3b).
25
26
27

28 29 *Key topics and approaches for analyses*

30
31 The EoI process and subsequent refinement resulted in a list of 10 research questions that
32 addressed determinants/correlates across the life course, had a balanced focus on physical activity,
33 sedentary behaviour or both, and required a variety of statistical approaches (e.g., using Bayesian
34 network modelling, Chi-squared automatic interaction detection analysis,[33] etc.). Work teams
35 around each research question were formed. A three-day writing retreat was organised in Ghent,
36 Belgium, to further define approaches, to progress specific questions as well as identify unresolved
37 issues regarding the pooling and harmonisation process.
38
39

40 41 *Reusability of target datasets*

42 The metadata contained within the compendium provide details on the potential level of reusability,
43 and what conditions must be met to facilitate access to the datasets (Figure 2). The compendium
44 provides metadata for all included datasets, regardless of the level of reusability of the data.
45

46
47 Of all 114 identified datasets, 43 were deemed to have potential in answering the 10 exemplar
48 questions posed and subsequently data access requests were made to the dataset owners. Of these,
49 actually, 12 datasets were not accessible during the timeframe of the DEDIPAC project, 10 dataset
50 owners did not respond to our request after a number of attempts, and 20 datasets (47%) agreed
51 that access was possible. At the time of analyses, data had been obtained from 14 dataset owners.
52
53
54
55
56
57
58
59
60

Table 3a: Breakdown of datasets reporting on physical activity correlates or determinants by stage of the life course (where applicable)

	Datasets including physical activity determinants (overall)	Biological	Psychological	Behavioural	Physical Environmental	Socio- cultural	Economic	Policy
Children only	8	6	2	6	2	6	1	3
Adolescents only	1	1	1	1	1	1	1	0
Adults only	11	7	7	3	3	1	1	0
Older Adults only	3	2	0	1	0	1	0	0
Children and Adolescents	5	2	3	3	1	0	0	0
All stages of life course	2	0	0	1	2	0	0	0
Children, Adults, Older Adults	3	2	2	2	1	0	0	0
Adolescents, Adults, Older Adults	3	1	0	1	2	2	1	0
Adults and Older Adults	11	7	4	5	2	2	1	0
Children and Adults	1	1	0	1	0	0	0	0
Children, Adolescents, Adults	3	2	3	1	2	2	1	0
Adolescents and Adults	2	0	1	1	0	0	0	0
Total	53	31	23	26	16	15	6	3

Table 3b: Breakdown of datasets reporting on sedentary correlates or determinants by stage of the life course (where applicable)

	Datasets including sedentary behaviour determinants (overall)	Biological	Psychological	Behavioural	Physical Environmental	Socio- cultural	Economic	Policy
Children only	5	1	0	4	0	1	1	0
Adults only	1	1	1	0	1	1	1	0
Children and Adolescents	3	0	0	3	0	0	0	0
All stages of life course	2	0	0	2	0	0	0	0
Children, Adults, Older Adults	2	2	1	0	0	0	1	0
Adolescents, Adults, Older Adults	1	1	0	0	1	1	1	0
Adults and Older Adults	2	1	1	0	1	0	0	0
Children, Adolescents, Adults	1	1	1	0	0	0	1	0
Total	17	7	4	9	3	3	5	0

Data Pooling and Harmonisation

The dataset owners that agreed to participate transferred the variables outlined in the data request form. The pooling of these data broadly followed the FAIR principles for interoperability, to facilitate the pooling of data across non-co-operating resources. The data was transferred in a compatible electronic format (most often SPSS, sometimes SAS or STATA). Data was then transferred to SPSS (IBM SPSS Statistics for Windows, Version 22.0. Armonk, NY, USA) and the original data was archived for backup purposes. The datasets were then distributed to those leading the analyses. Variables to be harmonised ranged from raw accelerometry data to summary data (e.g. proportion of the population above a defined threshold). As a proof of concept, the DEDIPAC DHP was used for one research question. This comprised the pooling of two large datasets and the first simple harmonisation steps encompassed the alignment of gender, education coding, etc.

Discussion

There is an increasing emphasis on the potential utility of existing datasets as a resource to answer research questions. While there are obvious benefits in terms of increased power for statistical analysis, this new and evolving approach has a number of challenges. The DEDIPAC consortium sought, for the first time, to examine the potential for data pooling and secondary data analysis in the measurement and determinants of dietary, physical activity and sedentary behaviours. In so doing, DEDIPAC has developed a valuable compendium of datasets that make it possible to ascertain the scope and scale of European projects in this research domain. Whereas low levels of physical activity and too much sitting are causing health problems worldwide, the focus of this paper is on the European setting. Many of the issues encountered regarding data pooling are expected to be similar for other regions, but the potential for data pooling and secondary data analyses may also partly be different for studies conducted in and by researchers from other regions, for example because of differences in rules and regulations regarding data sharing. Furthermore, the data pooling was conducted to further research into potential determinants of sedentary behaviour and physical activity, and such determinants themselves may be (partly) different between regions of the world, e.g. because of differences in affluence, infrastructure and cultural differences.

In our approach we applied the FAIR principles to provide guidance in the discovery and reuse of data for further investigation. We have identified more than one hundred potentially relevant European datasets through our outlined search approaches (*'Findability'*). It was possible to retrieve the metadata from these datasets on the types of variables, age groups under study, study design, measurement instruments used, time frame, etc. (*'Accessibility'*). These metadata were systematically detailed in the DEDIPAC compendium in order to be *'Interoperable'*. These datasets were, and are, being used to address 10 exemplar research questions (*'Reusability'*) via either direct analyses, de-centralised pooling and harmonisation of multiple datasets, or centralised pooling and harmonisation using a DHP. Therefore, we can conclude that, as a proof-of-concept, it is possible to apply the FAIR principles and successfully undertake research projects using existing data.

However, along the way, we encountered a number of significant challenges that were specific to the elements of the FAIR principles, but we also faced issues that applied to the pooling,

1
2
3 harmonising and/or analysing the secondary data. These issues should be carefully considered if
4 data pooling and harmonisation are to become a more central aspect of pan-European data analysis.
5

6 Firstly, data that does not exist is not possible to find. The inventory of data in the development of
7 the compendium brought to light that the number of existing datasets that may be relevant for
8 pooling, harmonising and secondary data analyses in this field of research was rather limited – or
9 very well hidden. There were especially few current datasets that contained a wide set of variables
10 (outcomes, independent variables, co-variables) to study behaviours and their underlying factors.
11 Therefore, the DEDIPAC project has highlighted a dearth of data on the determinants of physical
12 activity and sedentary behaviour and believe a pan-European cohort study with a focus on
13 behavioural determinants is required to address this deficiency.
14
15

16
17 Secondly, while we have demonstrated that it is feasible to retrieve metadata from datasets, it often
18 required repeated personal (e-mail) contact with dataset owners. In addition, and more importantly,
19 the reusability of datasets was limited to the small number of datasets to which access was granted.
20 In the timeframe of the DEDIPAC project it was challenging and time consuming to pursue access to
21 the datasets to address the 10 exemplar research questions and only 12% had provided access at the
22 time of writing. In the future, different approaches may be required to encourage dataset owners to
23 participate in pooled data analysis. The outcome of the DEDIPAC exemplar projects may provide the
24 evidence of the benefits and limitations that would provide dataset owners with a more solid basis
25 for engagement. It is also possible that funding agencies may mandate the sharing of publicly funded
26 data. Therefore, a continued and open discussion of the merits and limitations of data pooling is
27 necessary.
28
29
30

31 Thirdly, and beyond FAIR, the harmonization of core outcome measures under study (in most cases
32 sedentary behaviour and physical activity) was often problematic. Across the included and accessible
33 datasets, these outcomes were measured or operationalized in a variety of ways. The substantial
34 variation in assessment methods and operationalization of outcome variables across current
35 European studies (as illustrated in Table 2) not only hampered the practical harmonisation process,
36 but also presented comparability issues, as estimations of physical activity and sedentary behaviour
37 levels are known to differ based on the assessment method used.[21–24]
38
39

40
41 Fourthly, next to harmonisation issues of core outcomes (physical activity and/or sedentary
42 behaviour, in this context), our focus on determinants of physical activity and sedentary behaviours
43 meant that individual-level and contextual – more upstream – factors were to be taken into account.
44 It was possible to harmonise some of the core outcomes sometimes, and some key socio-
45 demographic variables (e.g. age, gender and educational background) could be harmonised. However,
46 harmonisation of the behavioural determinants was rarely possible and important co-variables even
47 less so. Thus, pooling often implied that certain variables could not be taken into account if they
48 were not measured across all included datasets, or could not be harmonised. In our opinion, the
49 different assessment methods impeded harmonisation.
50
51
52

53 Finally, the retrospective pooling and harmonization of variables requires a ‘flexible design’, since
54 very few established studies have used identical collection methods and procedures.[34] A flexible
55 design means that various categories of a variable such as, for instance, education attainment are
56 reduced to few, or even only two categories (e.g. higher/lower education). Hence, putting data from
57 respondents from different studies together in one dataset indeed increases the number of
58
59
60

1
2
3 observations and the absolute power of observations, but the promise that smaller associations can
4 be picked up is likely to be compromised. The added nuance sought for with pooling data may thus
5 be undone or even reversed by the – often rough – re-categorisation of variables during the
6 harmonisation process. This may be prevented in prospective harmonisation efforts – where
7 particular instruments, protocols and operationalisations are specified and aligned beforehand. Such
8 prospective harmonisation is especially required in surveillance systems to monitor regional
9 variations and temporal trends of health behaviours and health outcomes, although this is not often
10 done.[35, 36] In fact an important goal of DEDIPAC was to work towards harmonization of
11 measurement instruments to better enable and promote prospective data harmonization.[19]
12
13

14 15 16 17 **Conclusions**

18
19 The DEDIPAC project has identified and compiled a large number of pan-European projects related
20 to physical activity and sedentary behaviour. In the design and analysis of 10 exemplar projects,
21 using a variety of approaches, we identified a number of challenges including (i) the limited
22 availability of datasets that contain variables to examine the determinants of physical activity and
23 sedentary behaviour, (ii) the difficulty to establish communication with dataset owners or getting
24 their agreement to share data for analysis and (iii) a low harmonisation potential for the limited
25 number of variables that were available, especially for the potential determinants further upstream.
26 Compliance to FAIR data management and stewardship principles currently appears to be limited for
27 research in the field of physical activity and sedentary behaviour. It is recognized that researchers
28 should be facilitated, funded, requested, or required (e.g. by funding agencies) to share their data
29 and comply to the FAIR principles [18,37]. While recognising the importance of utilising existing data
30 it is equally, if not more important, to highlight the absence of data that would be needed to
31 investigate in detail the determinants of these behaviours. Not complying to the FAIR principles will
32 limit the reusability of relevant data. In turn, the lack of suitable data will severely limit the ability of
33 research to understand and predict these behaviours and inform policy. A bigger, targeted and more
34 standardized data collection is needed in order to maximize the potential of data pooling and
35 harmonisation. Until then, there are only narrow margins for determinant research to build and
36 harvest on previous work.
37
38
39
40
41
42
43

44 **Declarations**

45 *Competing interests*

46 The authors declare that they have no competing interests.
47

48 *Data sharing statement*

49 The data may be obtained from the authors for academic purposes
50
51

52 *Funding*

53
54 The preparation of this paper was supported by the DEterminants of Diet and Physical Activity
55 (DEDIPAC) knowledge hub. This work is supported by the Joint Programming Initiative 'Healthy Diet
56 for a Healthy Life'. The funding agencies supporting this work are: Belgium: Research Foundation –
57 Flanders; France: Institut National de la Recherche Agronomique (INRA); Germany: Federal Ministry
58
59
60

of Education and Research; Italy: Ministry of Education, University and Research (DEDIPAC F.S. 02.15.02 COD. B84G14000040008; CDR2.PRIN 2010/11 COD. 2010KL2Y73_003)/ Ministry of Agriculture Food and Forestry Policies; Ireland: The Health Research Board (HRB); The Netherlands: The Netherlands Organisation for Health Research and Development (ZonMw); The United Kingdom: The Medical Research Council (MRC).

Authors' contributions

JB, GC, SC, HvdP, JL, LC and CMD conceived the study. FL and AC assisted in data collection. AL and JK redesigned the DEDIPAC data warehouse. JL drafted the manuscript and AL, FL, MDC, HvdP, DOG, AC, LC, JK, JMO, SC, GC, JB and CMD contributed to iterative revisions.

References

- 1 WHO. Physical activity Factsheet N°385, Reviewed June 2016.
- 2 Wilmot EG, Edwardson CL, Achana FA, *et al.* Sedentary time in adults and the association with diabetes, cardiovascular disease and death: systematic review and meta-analysis. *Diabetologia* 2012;**56**:942-3.
- 3 Grontved A, Hu F. Television viewing and risk of type 2 diabetes, cardiovascular disease, and all-cause mortality: a meta-analysis. *JAMA* 2011;**305**:2448–55.
- 4 Biswas A, Oh PI, Faulkner GE, *et al.* Sedentary Time and Its Association With Risk for Disease Incidence, Mortality, and Hospitalization in Adults. *Ann Intern Med* 2015;**162**:123.
- 5 Lee IM, Shiroma EJ, Lobelo F, *et al.* Effect of physical inactivity on major non-communicable diseases worldwide: An analysis of burden of disease and life expectancy. *Lancet* 2012;**380**:219–29.
- 6 WHO, National Institute on Aging, National Institutes of Health. Global Health and Aging. 2011.
- 7 Bauman AE, Reis RS, Sallis JF, *et al.* Correlates of physical activity: why are some people physically active and others not? *Lancet* 2012;**380**:258–71.
- 8 Marques A, Martins J, Peralta M, *et al.* European adults' physical activity socio-demographic correlates: a cross-sectional study from the European Social Survey. *PeerJ* 2016;**4**:e2066.
- 9 Carlin A, Perchoux C, Puggina A, *et al.* A life course examination of the physical environmental determinants of physical activity behaviour: A 'Determinants of Diet and Physical Activity' (DEDIPAC) umbrella systematic literature review. *PLoS One* 2017;**12**:e0182083.
- 10 Condello G, Puggina A, Aleksovskaja K, *et al.* Behavioral determinants of physical activity across the life course: a 'Determinants of Diet and Physical Activity' (DEDIPAC) umbrella systematic literature review. *Int J Behav Nutr Phys Act* 2017;**14**:58.
- 11 Stierlin AS, De Lepeleere S, Cardon G, *et al.* A systematic review of determinants of sedentary behaviour in youth: a DEDIPAC-study. *Int J Behav Nutr Phys Act* 2015;**12**:133.

- 1
2
3 12 Chastin SFM, Buck C, Freiburger E, *et al.* Systematic literature review of determinants of
4 sedentary behaviour in older adults: a DEDIPAC study. *Int J Behav Nutr Phys Act* 2015;**12**:127.
5
6 13 O'Donoghue G, Perchoux C, Mensah K, *et al.* A systematic review of correlates of sedentary
7 behaviour in adults aged 18-65 years: a socio-ecological approach. *BMC Public Health*
8 2016;**16**:163.
9
10 14 Lakerveld J, Mackenbach JD. The upstream determinants of adult obesity. *Obes Facts*
11 2017;**10**:216–22.
12
13 15 Doiron D, Burton P, Marcon Y, *et al.* Data harmonization and federated analysis of
14 population-based studies: the BioSHaRE project. *Emerg Themes Epidemiol* 2013;**10**:12.
15
16 16 Pisani E, AbouZahr C. Sharing health data: good intentions are not enough. *Bull World Health*
17 *Organ* 2010;**88**:462–6.
18
19 17 Schofield PN, Eppig J, Huala E, *et al.* Sustaining the Data and Bioresource Commons. *Science*
20 2010;**330**:592–3.
21
22 18 Piwowar HA, Becich MJ, Bilofsky H, *et al.* Towards a data sharing culture: Recommendations
23 for leadership from academic health centers. *PLoS Med.* 2008;**5**:1315–9.
24
25 19 Lakerveld J, van der Ploeg HP, Kroeze W, *et al.* Towards the integration and development of a
26 cross-European research network and infrastructure: the DETERminants of Diet and Physical
27 ACTivity (DEDIPAC) Knowledge Hub. *Int J Behav Nutr Phys Act* 2014;**11**:143.
28
29 20 Loyen A, Van Der Ploeg HP, Bauman A, *et al.* European sitting championship: Prevalence and
30 correlates of self-reported sitting time in the 28 European Union Member States. *PLoS One*
31 2016;**11**: e0149320.
32
33 21 Loyen A, Verloigne M, Van Hecke L, *et al.* Variation in population levels of sedentary time in
34 European adults according to cross-European studies: a systematic literature review within
35 DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**:1–11.
36
37 22 Van Hecke L, Loyen A, Verloigne M, *et al.* Variation in population levels of physical activity in
38 European children and adolescents according to cross-European studies: a systematic
39 literature review within DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**:1–22.
40
41 23 Verloigne M, Loyen A, Van Hecke L, *et al.* Variation in population levels of sedentary time in
42 European children and adolescents according to cross-European studies: a systematic
43 literature review within DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**:1–30.
44
45 24 Loyen A, Van Hecke L, Verloigne M, *et al.* Variation in population levels of physical activity in
46 European adults according to cross-European studies: a systematic literature review within
47 DEDIPAC. *Int J Behav Nutr Phys Act* 2016;**13**:72.
48
49 25 Granda P, Blasczyk E. Data harmonization, guidelines for best practice in cross-cultural
50 surveys. *MI Surv Res Centre, Inst Soc Res Univ Michigan* 2010.
51
52 26 Gaye A, Marcon Y, Isaeva J, *et al.* DataSHIELD: Taking the analysis to the data, not the data to
53 the analysis. *Int J Epidemiol* 2014;**43**:1929–44.
54
55 27 Crosas M. The dataverse network®: An open-source application for sharing, discovering and
56 preserving data. *D-Lib Mag* 2011;**17**:1/2.
57
58 28 Sherar LB, Griew P, Esliger DW, *et al.* International children's accelerometry database (ICAD):
59
60

- 1
2
3 design and methods. *BMC Public Health* 2011;**11**:485.
4
5 29 Buffart LM, Kalter J, Chinapaw MJM, *et al.* Predicting Optimal Cancer Rehabilitation and
6 Supportive care (POLARIS): rationale and design for meta-analyses of individual patient data
7 of randomized controlled trials that evaluate the effect of physical activity and psychosocial
8 interventions on health. *Syst Rev* 2013;**2**:1–9.
9
10 30 Wilkinson MD, Dumontier M, Aalbersberg IJJ, *et al.* The FAIR Guiding Principles for scientific
11 data management and stewardship. *Sci Data* 2016;**3**:160018.
12
13 31 European Commission. *CORDIS* (Community Research and Development Information
14 Services).
15
16 32 Chastin SFM, De Craemer M, Lien N, *et al.* The SOS-framework (Systems of Sedentary
17 behaviours): an international consensus transdisciplinary framework and research priorities
18 for the study of determinants of sedentary behaviour and policy across the life course: a
19 DEDIPAC study. *Int J Behav Nutr Phys Act* 2016;**13**:83.
20
21 33 Lakerveld J, Loyen A, Schotman N, *et al.* Sitting too much: A hierarchy of socio-demographic
22 correlates. *Prev Med* 2017;**101**:77–83.
23
24 34 Doiron D, Parminder R, Ferretti V, *et al.* Facilitating collaborative research: Implementing a
25 platform supporting data harmonization and pooling. *Nor Epidemiol* 2012;**21**:221–4.
26
27 35 Bel-Serrat S, Huybrechts I, Thumann BF, *et al.* Inventory of surveillance systems assessing
28 dietary, physical activity and sedentary behaviours in Europe: a DEDIPAC study. *Eur J Public*
29 *Health* 2017;**56**:S25-32.
30
31 36 Brug J, van der Ploeg HP, Loyen A, *et al.* Towards better knowledge of the causes of the
32 causes: lessons learned from European Joint Programming in research on Determinants of
33 Diet and Physical Activity (DEDIPAC). Submitted
34
35 37 Genbank A. Let data speak to data. *Nature* 2005;**438**:531.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **Tables and Figures**
4
5
6

7 **Figure 1 capture:**
8

9 **Figure 1:** Schematic outline of the stepwise approach towards secondary data analyses
10
11

12
13 **Figure 2 capture:**
14

15 **Figure 2:** Flow diagram of accessibility to datasets
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

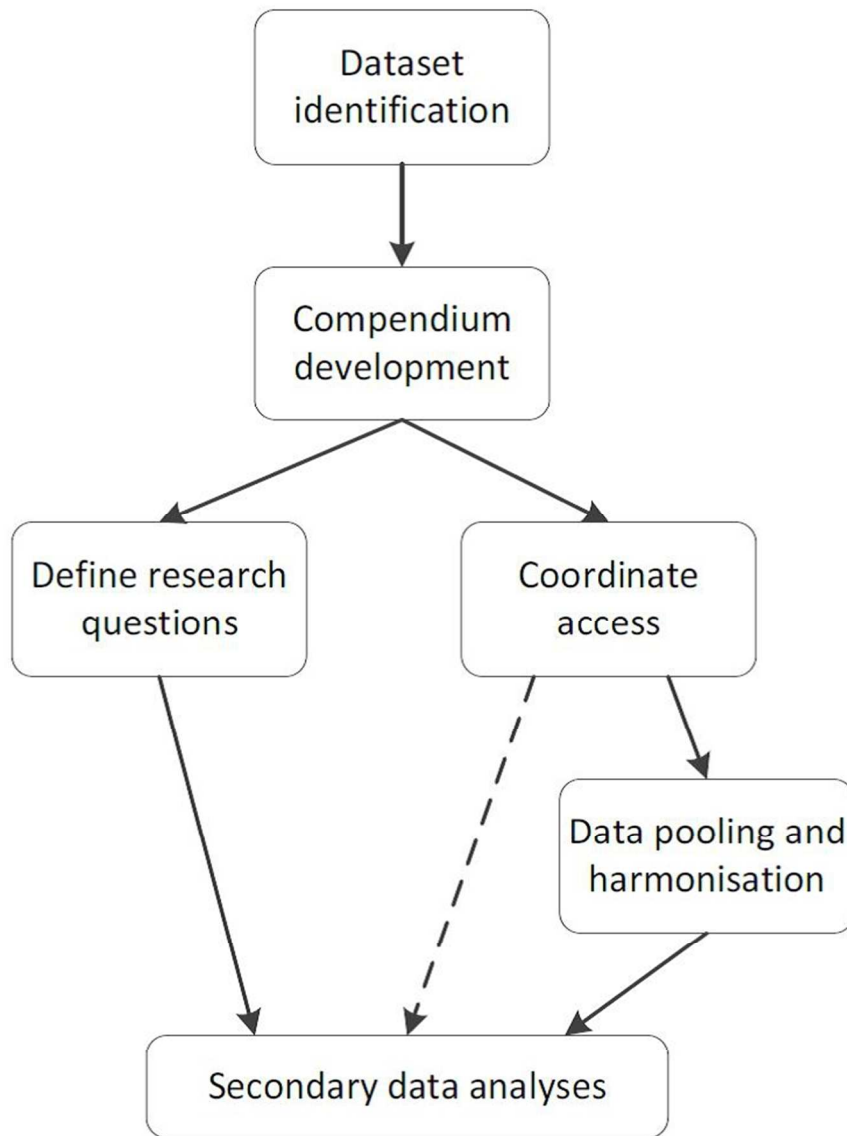
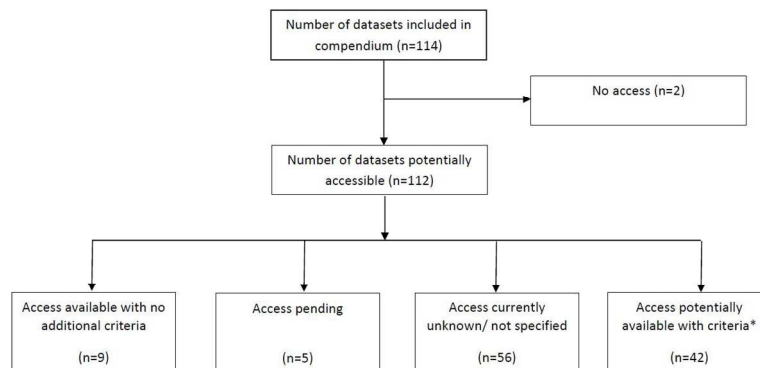


Figure 1: Schematic outline of the stepwise approach towards secondary data analyses

72x85mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



* Criteria applies regarding access to datasets, for example, Access with Permission, Access after Board Review and Written Material Transfer Agreement, Access to Raw Data at Local Site and/or Provision of Summary Tables for Meta-Analysis (note: for some datasets multiple criteria apply)

Figure 2: Flow diagram of accessibility to datasets

140x66mm (300 x 300 DPI)