#### Elsevier Editorial System(tm) for Vaccine Manuscript Draft

Manuscript Number:

Title: Quantifying outcome misclassification in multi-database studies: the case study of pertussis in the ADVANCE project

Article Type: SI: ADVANCE

Keywords: incidence of pertussis; event-finding algorithms; event misclassification; positive predictive value

Corresponding Author: Dr. Rosa Gini, Ph.D.

Corresponding Author's Institution: Agenzia regionale di sanità della Toscana

First Author: Rosa Gini, Ph.D.

Order of Authors: Rosa Gini, Ph.D.; Caitlin N Dood; Kaatje Bollaerts; Claudia Bartolini; Giuseppe Roberto; Claudia Huerta-Alvarez; Elisa Martín-Merino; Talita Duarte-Salles; Gino Picelli; Lara Tramontan; Giorgia Danieli; Ana Correa; Chris McGee; Benedikt F Becker; Charlotte Switzer; Sonja Gandhi-Banga; Jorgen Bauwens; Nicoline A van der Maas; Gianfranco Spiteri; Emmanouela Sdona; Daniel Weibel; Miriam Sturkenboom

Abstract: Background

The Accelerated Development of VAccine beNefit-risk Collaboration in Europe (ADVANCE) is a public-private collaboration aiming to develop and test a system for rapid benefit-risk (B/R) monitoring of vaccines using European healthcare databases. Event misclassification can result in biased estimates and contribute to heterogeneity in results. Here we report the impact of different event-finding algorithms for Bordetella pertussis (BorPer) on the estimated incidence rates (IRs) and algorithm validity.

Methods

Four participating databases retrieved data from primary care (PC) setting: (BIFAP: Spain), THIN and RCGP RSC: UK) and PEDIANET: Italy); the fifth SIDIAP (Spain) from both PC and hospital settings. The algorithms were defined by setting, data domain (diagnoses, drugs, or tests) and concept sets (specific or unspecified pertussis). BorPer IRs were estimated in children aged 0-14 years enrolled in 2012 and 2014 and followed up until the end of each calendar year and compared with IRs of confirmed pertussis from the ECDC surveillance system (TESSy). Results

The number of cases and the estimated BorPer IRs per 100,000 person-years in PC, using data representing 3,173,268 person-years, were 0 (IR=0.0), 21 (IR=4.3), 21 (IR=5.1), 79 (IR=5.7), and 2 (IR=2.3) in BIFAP, SIDIAP, THIN, RCGP RSC and PEDIANET respectively. The IRs for combined specific/unspecified pertussis were higher and were comparable with data from TESSy, except PEDIANET. In SIDIAP the estimated IR was 45.0 when discharge diagnoses were included. The sensitivity and positive predictive value of combined PC specific and unspecific diagnoses for BorPer cases in SIDIAP were 85% and 72%, respectively, based on overlap between hospital and PC diagnoses (adjusted IR=35.5). Conclusion This study demonstrated the value of quantifying the impact of different event-finding algorithms across databases and the possibility of benchmarking with disease surveillance data as well as assessing validity estimates when data from different settings can be linked. Dr Gregory A Poland Editor-in-Chief, Vaccine

18 January 2019

Dear Dr Poland

We are pleased to submit our paper '*Quantifying outcome misclassification in multi-database studies: the case study of pertussis in the ADVANCE project*' to your Journal, Vaccine, for the ADVANCE supplement. This paper describes the analyses we have performed to quantify the impact of different event-finding algorithms across databases on outcome misclassification. It is the 8th of the series of 10 papers that will be included in the supplement.

On behalf of all co-authors

Dr Rosa Gini

I, Dr. Rosa Gini, declare that all authors have seen and approved the final version of the manuscript being submitted. We warrant that the article is our original work that has not been previously published and is not under consideration for publication elsewhere

Dr Rosa Gini

Name	Institute	email
Gianluca Trifirò	Università di Messina, Italy	trifirog@unime.it
Jeff Brown	Department of Population Medicine (DPM) at Harvard Medical School and the Harvard Pilgrim Health Care Institute, USA	jeff_brown@harvardpilgrim.org
Gillian Hall	Gillian Hall Epidemiology LTD, UK	gillian.hall@gchall.com
Patrick Ryan	Observational Health Data Sciences and Informatics (OHDSI)	ryan@ohdsi.org

#### **Declaration of interests**

□ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

⊠ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Rosa Gini, Caitlin Dodd, Kaatje Bollaerts, Claudia Bartolini, Giuseppe Roberto, Consuelo Huerta-Alvarez, Elisa Martín-Merino, Talita Duarte-Salles, Gino Picelli, Lara Tramontan, Giorgia Danieli, Ana Correa, Chris McGee, Benedikt Becker, Charlotte Switzer, Jorgen Bauwens, Nicoline van der Maas, Gianfranco Spiteri, Emmanouela Sdona declared no conflicts of interest. Sonja Gandhi-Banga declared that she works for Sanofi Pasteur and holds company shares. Daniel Weibel declared that he has received personal fees from GSK for work unrelated to the submitted work. Miriam Sturkenboom declared that she has received grants from Novartis, CDC and Bill & Melinda Gates Foundation for work unrelated to the submitted work.

#### 1 Abstract

2 Background

3 The Accelerated Development of VAccine beNefit-risk Collaboration in Europe

4 (ADVANCE) is a public-private collaboration aiming to develop and test a system for rapid

5 benefit-risk (B/R) monitoring of vaccines using European healthcare databases. Event

6 misclassification can result in biased estimates and contribute to heterogeneity in results. Here

7 we report the impact of different event-finding algorithms for *Bordetella pertussis* (BorPer)

8 on the estimated incidence rates (IRs) and algorithm validity.

9 *Methods* 

10 Four participating databases retrieved data from primary care (PC) setting: (BIFAP: Spain),

11 THIN and RCGP RSC: UK) and PEDIANET: Italy); the fifth SIDIAP (Spain) from both PC

12 and hospital settings. The algorithms were defined by setting, data domain (diagnoses, drugs,

13 or tests) and concept sets (specific or unspecified pertussis). BorPer IRs were estimated in

14 children aged 0-14 years enrolled in 2012 and 2014 and followed up until the end of each

15 calendar year and compared with IRs of confirmed pertussis from the ECDC surveillance

16 system (TESSy).

17 Results

18 The number of cases and the estimated BorPer IRs per 100,000 person-years in PC, using data

19 representing 3,173,268 person-years, were 0 (IR=0.0), 21 (IR=4.3), 21 (IR=5.1), 79 (IR=5.7),

20 and 2 (IR=2.3) in BIFAP, SIDIAP, THIN, RCGP RSC and PEDIANET respectively. The IRs

21 for combined specific/unspecified pertussis were higher and were comparable with data from

22 TESSy, except PEDIANET. In SIDIAP the estimated IR was 45.0 when discharge diagnoses

23 were included. The sensitivity and positive predictive value of combined PC specific and

unspecific diagnoses for BorPer cases in SIDIAP were 85% and 72%, respectively, based on

25 overlap between hospital and PC diagnoses (adjusted IR=35.5).

# 26 Conclusion

- 27 This study demonstrated the value of quantifying the impact of different event-finding
- 28 algorithms across databases and the possibility of benchmarking with disease surveillance
- 29 data as well as assessing validity estimates when data from different settings can be linked.

1	Quantifying outcome misclassification in multi-database studies: the case study of
2	pertussis in the ADVANCE project
3	Rosa <b>Gini</b> <sup>a</sup> , Caitlin N <b>Dodd</b> <sup>b,c</sup> , Kaatje <b>Bollaerts</b> <sup>d</sup> , Claudia <b>Bartolini</b> <sup>a</sup> , Giuseppe <b>Roberto</b> <sup>a</sup> ,
4	Consuelo Huerta-Alvarez <sup>e</sup> , Elisa Martín-Merino <sup>e</sup> , Talita Duarte-Salles <sup>f</sup> , Gino Picelli <sup>g</sup> , Lara
5	<b>Tramontan</b> <sup>g,h</sup> , Giorgia <b>Danieli</b> <sup>g,h</sup> , Ana <b>Correa</b> <sup>i</sup> , Chris <b>McGee</b> <sup>i,j</sup> , Benedikt F H <b>Becker</b> <sup>b</sup> ,
6	Charlotte Switzer <sup>k1</sup> , Sonja Gandhi-Banga <sup>k</sup> , Jorgen Bauwens <sup>1,m,n</sup> , Nicoline A T van der
7	Maas <sup>m,n</sup> , Gianfranco Spiteri <sup>o</sup> , Emmanouela Sdona <sup>o2</sup> , Daniel Weibel <sup>b</sup> , Miriam
8	<b>Sturkenboom</b> <sup>c,d,p</sup>
9	<sup>a</sup> Agenzia regionale di sanità della Toscana, Osservatorio di epidemiologia, Florence, Italy
10	(rosa.gini@ars.toscana.it; claudia.bartolini@ars.toscana.it; giuseppe.roberto@ars.toscana.it)
11	<sup>b</sup> Erasmus University Medical Center, Post box 2040, 3000 CA Rotterdam, Netherlands
12	(caitlinndodd@gmail.com; benedikt.becker@posteo.de; d.weibel@erasmusmc.nl)
13	<sup>c</sup> Julius Global Health, University Medical Center, Utrecht, Heidelberglaan 100, The
14	Netherlands (caitlinndodd@gmail.com; m.c.j.sturkenboom@umcutrecht.nl)
15	<sup>d</sup> P95 Epidemiology and Pharmacovigilance, Koning Leopold III laan 1 3001 Heverlee,
16	Belgium (kaatje.bollaerts@p-95.com; m.c.j.sturkenboom@umcutrecht.nl)
17	<sup>e</sup> BIFAP database, Spanish Agency of Medicines and Medical Devices, Madrid, Spain
18	(chuerta@aemps.es; emartinm@aemps.es)
19	<sup>f</sup> Institut Universitari d'Investigació en Atenció Primària Jordi Gol (IDIAP Jordi Gol),
20	Barcelona, Spain (tduarte@idiapjgol.org)
21	<sup>g</sup> Epidemiological Information for Clinical Research from an Italian Network of Family
22	Paediatricians (PEDIANET), Padova, Italy (g.picelli@virgilio.it;
23	ltramontan@consorzioarsenal.it; gdanieliconsorzioarsenal@gmail.com)

 <sup>&</sup>lt;sup>1</sup> Current address: McMaster University, 1280 Main St W, Hamilton, ON L8S 4L8, Canada.
 <sup>2</sup> Current address: Karolinska Institutet, Solnavägen 1, 171 77 Solna, Sweden

- <sup>h</sup>Consorzio Arsenal.IT, Veneto Region, Italy (ltramontan@consorzioarsenal.it;
- 25 gdanieliconsorzioarsenal@gmail.com)
- <sup>1</sup>University of Surrey, Guildford, Surrey GU2 7XH, UK (s.lusignan@surrey.ac.uk;
- 27 c.mcgee@surrey.ac.uk)
- <sup>j</sup>Royal College of General Practitioners, Research and Surveillance Centre, 30 Euston Square,
- 29 London NW1 2FB, UK (accorrea1@googlemail.com; c.mcgee@surrey.ac.uk)
- <sup>k</sup>Sanofi Pasteur, 1755 Steeles Ave W, North York, ON M2R 3T4, Canada
- 31 (switzer.charlotte@gmail.com; sonja.banga@sanofi.com)
- <sup>32</sup> <sup>1</sup>University Children's Hospital, Basel, Switzerland (j.bauwens@brightoncollaboration.org)
- <sup>33</sup> <sup>m</sup>University of Basel, Switzerland (j.bauwens@brightoncollaboration.org;
- 34 nicoline.van.der.maas@rivm.nl)
- <sup>n</sup>Brighton Collaboration Foundation, Switzerland (j.bauwens@brightoncollaboration.org;
- 36 nicoline.van.der.maas@rivm.nl)
- <sup>o</sup>European Centre for Disease Prevention and Control, Gustav III's Boulevard 40, 16973
- 38 Solna, Sweden (Gianfranco.Spiteri@ecdc.europa.eu; emysdona@gmail.com)
- 39 <sup>P</sup>VACCINE.GRID Foundation, Spitalstrasse 33, Basel, Switzerland
- 40 (m.c.j.sturkenboom@umcutrecht.nl)

1	Quantifying outcome misclassification in multi-database studies: the case study of
2	pertussis in the ADVANCE project
3	Rosa Gini <sup>a</sup> , Caitlin N Dodd <sup>b,c</sup> , Kaatje Bollaerts <sup>d</sup> , Claudia Bartolini <sup>a</sup> , Giuseppe Roberto <sup>a</sup> ,
4	Consuelo Huerta-Alvarez <sup>e</sup> , Elisa Martín-Merino <sup>e</sup> , Talita Duarte-Salles <sup>f</sup> , Gino Picelli <sup>g</sup> , Lara
5	<b>Tramontan</b> <sup>g,h</sup> , Giorgia <b>Danieli</b> <sup>g,h</sup> , Ana <b>Correa</b> <sup>i</sup> , Chris <b>McGee</b> <sup>i,j</sup> , Benedikt F H <b>Becker</b> <sup>b</sup> ,
6	Charlotte Switzer <sup>k1</sup> , Sonja Gandhi-Banga <sup>k</sup> , Jorgen Bauwens <sup>l,m,n</sup> , Nicoline A T van der
7	<b>Maas</b> <sup>m,n</sup> , Gianfranco <b>Spiteri</b> <sup>o</sup> , Emmanouela <b>Sdona</b> <sup>o2</sup> , Daniel <b>Weibel</b> <sup>b</sup> , Miriam
8	<b>Sturkenboom</b> <sup>c,d,p</sup>
9	<sup>a</sup> Agenzia regionale di sanità della Toscana, Osservatorio di epidemiologia, Florence, Italy
10	(rosa.gini@ars.toscana.it; claudia.bartolini@ars.toscana.it; giuseppe.roberto@ars.toscana.it)
11	<sup>b</sup> Erasmus University Medical Center, Post box 2040, 3000 CA Rotterdam, Netherlands
12	(caitlinndodd@gmail.com; benedikt.becker@posteo.de; d.weibel@erasmusmc.nl)
13	<sup>c</sup> Julius Global Health, University Medical Center, Utrecht, Heidelberglaan 100, The
14	Netherlands (caitlinndodd@gmail.com; m.c.j.sturkenboom@umcutrecht.nl)
15	<sup>d</sup> P95 Epidemiology and Pharmacovigilance, Koning Leopold III laan 1 3001 Heverlee,
16	Belgium (kaatje.bollaerts@p-95.com; m.c.j.sturkenboom@umcutrecht.nl)
17	<sup>e</sup> BIFAP database, Spanish Agency of Medicines and Medical Devices, Madrid, Spain
18	(chuerta@aemps.es; emartinm@aemps.es)
19	<sup>f</sup> Institut Universitari d'Investigació en Atenció Primària Jordi Gol (IDIAP Jordi Gol),
20	Barcelona, Spain (tduarte@idiapjgol.org)
21	<sup>g</sup> Epidemiological Information for Clinical Research from an Italian Network of Family
22	Paediatricians (PEDIANET), Padova, Italy (g.picelli@virgilio.it;
23	ltramontan@consorzioarsenal.it; gdanieliconsorzioarsenal@gmail.com)

 <sup>&</sup>lt;sup>1</sup> Current address: McMaster University, 1280 Main St W, Hamilton, ON L8S 4L8, Canada.
 <sup>2</sup> Current address: Karolinska Institutet, Solnavägen 1, 171 77 Solna, Sweden

- <sup>h</sup>Consorzio Arsenal.IT, Veneto Region, Italy (ltramontan@consorzioarsenal.it;
- 25 gdanieliconsorzioarsenal@gmail.com)
- <sup>1</sup>University of Surrey, Guildford, Surrey GU2 7XH, UK (s.lusignan@surrey.ac.uk;
- 27 c.mcgee@surrey.ac.uk)
- <sup>j</sup>Royal College of General Practitioners, Research and Surveillance Centre, 30 Euston Square,
- 29 London NW1 2FB, UK (accorrea1@googlemail.com; c.mcgee@surrey.ac.uk)
- <sup>k</sup>Sanofi Pasteur, 1755 Steeles Ave W, North York, ON M2R 3T4, Canada
- 31 (switzer.charlotte@gmail.com; sonja.banga@sanofi.com)
- <sup>32</sup> <sup>1</sup>University Children's Hospital, Basel, Switzerland (j.bauwens@brightoncollaboration.org)
- <sup>33</sup> <sup>m</sup>University of Basel, Switzerland (j.bauwens@brightoncollaboration.org;
- 34 nicoline.van.der.maas@rivm.nl)
- <sup>n</sup>Brighton Collaboration Foundation, Switzerland (j.bauwens@brightoncollaboration.org;
- 36 nicoline.van.der.maas@rivm.nl)
- <sup>o</sup>European Centre for Disease Prevention and Control, Gustav III's Boulevard 40, 16973
- 38 Solna, Sweden (Gianfranco.Spiteri@ecdc.europa.eu; emysdona@gmail.com)
- 39 <sup>P</sup>VACCINE.GRID Foundation, Spitalstrasse 33, Basel, Switzerland
- 40 (m.c.j.sturkenboom@umcutrecht.nl)
- 41
- 42

## 43 Abbreviations used

- 44 BorPer: Bordetella pertussis
- 45 B/R: benefit-risk
- 46 ECDC: European Centre for Disease Prevention and Control
- 47 IR: incidence rate
- 48 PC: primary care
- 49 PPV: positive predictive value
- 50

## 51 Abstract

#### 52 Background

53 The Accelerated Development of VAccine beNefit-risk Collaboration in Europe

54 (ADVANCE) is a public-private collaboration aiming to develop and test a system for rapid

55 benefit-risk (B/R) monitoring of vaccines using European healthcare databases. Event

56 misclassification can result in biased estimates and contribute to heterogeneity in results. Here

57 we report the impact of different event-finding algorithms for *Bordetella pertussis* (BorPer)

58 on the estimated incidence rates (IRs) and algorithm validity.

59 *Methods* 

60 Four participating databases retrieved data from primary care (PC) setting: (BIFAP: Spain),

61 THIN and RCGP RSC: UK) and PEDIANET: Italy); the fifth SIDIAP (Spain) from both PC

62 and hospital settings. The algorithms were defined by setting, data domain (diagnoses, drugs,

or tests) and concept sets (specific or unspecified pertussis). BorPer IRs were estimated in

64 children aged 0-14 years enrolled in 2012 and 2014 and followed up until the end of each

65 calendar year and compared with IRs of confirmed pertussis from the ECDC surveillance

66 system (TESSy).

67 Results

68 The number of cases and the estimated BorPer IRs per 100,000 person-years in PC, using data 69 representing 3,173,268 person-years, were 0 (IR=0.0), 21 (IR=4.3), 21 (IR=5.1), 79 (IR=5.7), 70 and 2 (IR=2.3) in BIFAP, SIDIAP, THIN, RCGP RSC and PEDIANET respectively. The IRs 71 for combined specific/unspecified pertussis were higher and were comparable with data from 72 TESSy, except PEDIANET. In SIDIAP the estimated IR was 45.0 when discharge diagnoses 73 were included. The sensitivity and positive predictive value of combined PC specific and unspecific diagnoses for BorPer cases in SIDIAP were 85% and 72%, respectively, based on 74 75 overlap between hospital and PC diagnoses (adjusted IR=35.5).

# 76 Conclusion

- 77 This study demonstrated the value of quantifying the impact of different event-finding
- algorithms across databases and the possibility of benchmarking with disease surveillance
- 79 data as well as assessing validity estimates when data from different settings can be linked.
- 80
- 81 Keywords: incidence of pertussis; event-finding algorithms; event misclassification; positive
  82 predictive value
- 83

## 84 **1. Introduction**

85 ADVANCE is a public-private collaboration aiming to develop and test a system for rapid benefit-risk (B/R) monitoring of vaccines using existing healthcare databases in Europe [1] 86 87 (see Appendix for list of consortium members). These databases have proven very useful for 88 studying drug effects and are commonly used in pharmacoepidemiology [2]. 89 Identifying events, such as vaccine-preventable diseases, adverse events of interest, co-90 morbidities and exposure to vaccination, is a pivotal first step in vaccine B/R studies. Since 91 there is limited or no control over the primary data collection when using existing healthcare 92 databases, event retrieval is usually not perfect. Individuals who experienced the event might 93 not be retrieved, for example if an individual is admitted to hospital for the event but no 94 primary care (PC) diagnosis is recorded, the event will not be retrieved from PC databases 95 and some individuals might be identified as having the event when in fact they did not. In a 96 PC database, this typically happens when the physician had only a suspicion, or if it was a 97 ruled-out diagnosis.

98 Researchers who access these databases usually develop their own methods to identify events 99 of interest, which are not always fully transparent [3, 4]. Events may be retrieved by 100 combining information from different settings (e.g., PC and hospital) and data domains, for 101 example diagnostic codes, drugs as proxies (e.g. in the case of diabetes), or laboratory 102 measurements. Use of information from more than one data domain, compared with using 103 diagnoses information only, can alter the sensitivity and positive predictive value (PPV) of the 104 event-finding algorithm. This alteration may happen differently in different databases, due to 105 the local characteristics of the database, the population, or the healthcare system. 106 It is well established that misclassification of events (false positives or false negatives) can 107 introduce bias in epidemiological studies, which can be corrected, to some extent, using 108 statistical methods [5-7]. However, to correct this bias, some validity parameters such as

sensitivity and PPV are required [8]. For this a gold standard and chart reviews are required,which generally make it costly and time-consuming.

111 In an attempt to develop a systematic approach to quantifying the impact of using different 112 event identification algorithms in multi-national, multi-database studies, the *component* 113 algorithm strategy was developed (Roberto 2016): a set of standardized algorithms, called 114 *components*, are defined and applied in each database. The impact of different algorithms on 115 the resulting estimates of disease occurrence is subsequently measured [9]. In this study we 116 aimed to refine this strategy by further standardizing the process, by developing and applying 117 novel formulas, by using benchmark data from another source and by using a data source 118 which had data from two settings. Since the proof-of-concept studies of ADVANCE focused 119 on pertussis, we used this event as case study.

#### 120 **2. Methods**

#### 121 **2.1.** Bordetella pertussis disease information

122 Bordetella pertussis causes pertussis, a vaccine-preventable infectious disease of the 123 respiratory tract. Symptoms include paroxysms of cough typically lasting from 1 to 6 weeks 124 or more and these may be milder in adolescents or immunised children [10, 11]. Several tests 125 are available to confirm Bordetella pertussis infection, including culture (which takes up to 14 126 days), serology and nucleic acid amplification tests. Pertussis is a notifiable infectious disease 127 and cases should be reported to the national surveillance system in all the countries involved 128 in ADVANCE. European Union member states are required to report available data on 129 pertussis cases to the European Centre for Disease Prevention and Control (ECDC). A 130 standardised case definition is used which classifies cases based on clinical, epidemiological 131 and laboratory criteria [12]. All national reports are submitted to the European Surveillance 132 System database (TESSy) managed by the ECDC [13].

### 133 **2.2.** Data sources

134 We assessed the impact of different event-finding algorithms using five databases that

135 participated in the ADVANCE proof-of-concept studies: BIFAP and SIDIAP (Spain),

136 PEDIANET (Italy) and RCGP RSC and THIN (United Kingdom). All databases were

137 population-based with data from electronic medical records in the PC setting. In SIDIAP, the

138 analyses were restricted to the population in this PC database that could be linked to hospital

139 discharge records. Surveillance data on pertussis were obtained from the TESSy surveillance

140 system through ECDC, a partner of ADVANCE.

## 141 2.3. Study population and study design

142 We used a dynamic cohort study design to study the impact of different event-finding 143 algorithms on the estimated pertussis IRs. Due to the methodological nature of this study, to 144 enable us to explore in more detail a number of strategies, we included a larger cohort in the 145 study population than that in the other ADVANCE studies. Therefore, children aged 0 to 14 years who were registered in the participating databases entered the study cohorts on 1<sup>st</sup> 146 January 2012 and 1<sup>st</sup> January 2014, and were followed up during 2012 and 2014, respectively. 147 148 Children who were born during 2012 or 2014 were followed up from birth until the end of the 149 calendar year. Children who were older than 14 years at any point in the follow-up were 150 excluded. To exclude any previous cases that had been notified before the start of the study 151 period, children who had a record of one of the components of pertussis during the two years 152 prior to one of the cohort entry dates were excluded, unless the component referred to the data 153 domain of drugs (see below for more details on the component definition).

154 2.4. Selection of component algorithms

155 A component algorithm is a standardised event-finding algorithm specified by three 156 characteristics: the *setting* of primary data collection (PC or hospital), the *data domain* 157 involved in the algorithm, and the *set of concepts* used to find the codes used to query the

158 database [9]. The sets of concepts were created by aggregating the codes that were obtained 159 from an initial proposed list, completed with a literature review and pertussis case definitions 160 [2, 13]. The CodeMapper tool was used to support the process [14]. Labelling and 161 classification of identified concepts, as well as the construction of the components, were 162 conducted by one of the authors who is a pertussis expert (NvdM). As a result, seven concept 163 sets were created (Table 1) [15, 16]. In particular, two sets of concepts belonged to the 164 diagnoses data domain: the set labelled '(Bordetella pertussis)' included three concepts which 165 specifically indicated Bordetella pertussis as the causative agent of the infection, while the set 166 labelled '(pertussis unspecified)' included five concepts indicating unspecified pertussis. The 167 corresponding codes and free text keywords are given in **Supplementary Table 1**.

The primary components associating concepts with settings (PC and hospital) are described in **Table 2.** Some secondary components, combining primary components in pre-defined temporal relations (e.g., symptoms in the presence of a drug prescription in the previous 30 days) were also created.

172 **2.5.** Analysis

173 Each database manager received an R-coded programme (quality checked by double-coding 174 against Stata) which was programmed using the pre-specified common data model [1]. These 175 programmes produced aggregated outputs, which were then transferred to the remote research 176 environment. Event-finding algorithms were created as logical combinations of individual 177 components using Boolean operators. For example, the two components 'PC diagnosis, 178 specific' and 'PC diagnosis, unspecified' were combined in one component: 'PC specific OR 179 unspecified diagnoses', which detected all individuals that were positive for either of the 180 original components. Based on the different event-finding algorithms, incidence rates (IRs) 181 were estimated using the number of persons retrieved with the respective events as numerator 182 and the follow-up person-time as denominator (see Supplementary File 1).

183 Age and country-specific incidences per 100,000 person-years of confirmed BorPer for both 184 2012 and 2014 were calculated for children aged 0-14 years. The calculations used the 185 reported confirmed cases in the TESSy surveillance system in 2012 and 2014 as the 186 numerator, and person-time from population distributions in EUROSTAT for 2012 and 2014 187 as the denominator [17]. Exact Poisson confidence intervals (95% CI) were calculated [18]. 188 Some formulae link the true proportion of BorP and/or validity indices with each other and 189 with the observed proportion of the component algorithms (**Table 3**). These formulas are 190 explained in Supplementary File 2.

In this study we considered  $\Pi = IR$  (see **Supplementary File 1**) and we assumed that for all algorithms A and B, the proportion of true positives among those detected by both algorithms (PPV of A **AND** B), was the same as the PPV of A OR B, whichever was the highest, which may be considered the most conservative assumption.

195 Since the concept set labelled 'Bordetella pertussis' was composed of codes explicitly 196 mentioning the bacterium, we considered that components based on this had a high likelihood 197 of extracting true cases. Therefore we considered it was conservative to assume that the PPV 198 for 'PC diagnosis, specific' and for 'inpatient diagnosis, specific' was 90%. We explored two scenarios for the cases extracted by the components associated with the concept set labelled 199 200 'pertussis unspecified', assuming that the PPV was 70% or 50%. The PPV for the component 201 'positive laboratory results' was assumed to be 100%. Finally, we assumed that all true cases 202 in SIDIAP were recorded in at least one of the diagnosis or laboratory-based components: this 203 assumption may overestimate sensitivity. Based on this and on the formulae in Table 3, we 204 derived sensitivity and PPV estimates for the algorithm 'PC specific OR unspecified 205 diagnosis' in SIDIAP, and the adjusted IR of BorPer in the study population.

206 **3. Results** 

#### 207 3.1. Study population

- 208 We followed 3,173,268 person-years of children during the study period: 488,847 from the
- SIDIAP database, 796,324 from BIFAP, 88,754 from PEDIANET, 1,387,939 from THIN and
- 210 411,404 from RCGP RSC (**Table 4**). The percentages of children aged 0 or 1 years in the
- 211 population aged 0-14 years in Spain were 12.1% and 16.1% in SIDIAP in BIFAP,
- respectively, compared with 13.5% in the EUROSTAT Spanish population. In the UK the
- 213 percentages were 15.1% and 14.8% in RCGP RSC and THIN 13.0%, respectively, compared
- with 14.3% in the EUROSTAT UK population and in PEDIANET (vs 12.9%); in (vs 14.3%).

## 215 3.2. Incidence rates estimated by the algorithms

216 The IRs for the component and composite algorithms, as well as the benchmark IRs from the

- 217 TESSy surveillance system are illustrated in Figure 1 and documented in **Table 4**. The IRs
- estimated from the TESSy surveillance system in 2012 and 2014 for children aged 0-14 years
- 219 were 21.2 (95% CI: 20.5; 22.0) for Spain, 13.4 (95% CI: 13.0; 13.9) for the United Kingdom,
- and 5.4 (95% CI: 5.1; .8) for Italy. The number of cases of 'PC diagnosis, specific' (and IRs
- 221 per 100,000 PY) were 0 (0.0), 21 (4.3), 21 (5.1), 79 (5.7), and 2 (2.3) in the BIFAP, SIDIAP,
- 222 RCGP RSC, THIN and PEDIANET databases, respectively. The component 'PC diagnosis,
- 223 unspecified' had a higher IR in all databases, and combining the two components (one OR the
- other) increased the number of cases detected and the IRs to 135 (IR=17.0), 194 (IR=39.6), 39
- (IR=43.9), 246 (IR=17.7), and 91 (IR=22.1), respectively. In BIFAP, SIDIAP, RCGP RSC
- and THIN, when taking into account that the unspecified component may have captured some
- false positives, the IRs were comparable with the corresponding IR from the TESSy
- 228 surveillance system (17.0 vs 21.2; 39.6 vs. 21.2; 22.1 vs. 13.4; 17.7 vs. 13.4, respectively). In
- 229 PEDIANET the composite IR was much higher than the IR from the TESSy database (43.9 vs
- 230 5.4).

231 SIDIAP was the only database in which data from both the PC and hospital settings could be

232 linked. The total number of cases in 'PC OR inpatient diagnosis' in SIDIAP was 220

233 (IR=45.0), including 26 (12%) that had not been identified in the PC setting. Unlike in the PC

setting, where most of the diagnoses were unspecified, in the inpatient setting there were

around half specific and half unspecified diagnoses.

In BIFAP, the 'symptoms and drugs within 30 days' component identified 122 cases with an

237 IR of 15.3 per 100,000 PYs. When this component was combined with PC diagnoses, the IR

increased to 32.0, which was higher than the reference IR which was 21.2. Almost none of the

children aged 0 or 1 year old in 'symptoms in infants' in any database had a corresponding

240 prescription or dispensing of macrolides in the 'symptoms in infants and drugs within 30

241 days' component.

242 The 'test' component was available in all databases and had a relatively high IR (from 4.8 in

243 BIFAP to 42.8 in PEDIANET). 'Positive laboratory results' were only available in SIDIAP

and THIN, with only 19 and 3 cases, respectively. In SIDIAP, 3 of the 19 cases were notcaptured by a diagnosis in either primary care or hospital settings.

In **Supplementary Figure 1** and **Supplementary Table 2**, the analysis was repeated for infants (children aged 0 or 1). The IRs in this subpopulation were around three times higher than the IRs in the overall study population. The findings confirmed the relationship between components observed in the general study population, with the exection of 'PC OR inpatient diagnosis' in SIDIAP (n=98), where 25.5% (n=25) were not retrieved from the PC setting, vs 11.8% in the overall study population.

252 In SIDIAP we explored two scenarios, corresponding to different assumptions for PPV of 'PC

diagnosis, unspecified' and of 'inpatient diagnosis, unspecified': in the first, this was 70%, in

the second 50%. As a consequence, in the first scenario 'PC specific OR unspecified

diagnosis' had a PPV of 72% (or, in the second: 54%) and a sensitivity of 85% (or, in the

second: 83%). Based on this estimate, the adjusted IR of BorPer in the SIDIAP study
population was 35.5 per 100,000 PY (or, in the second scenario: 25.9) vs the TESSy
surveillance system IR 21.2.

259 **4. Discussion** 

260 We assessed several algorithms as potential strategies to detect cases of pertussis and thus 261 estimate the IR in five European healthcare databases. The IRs estimated by these algorithms 262 were heterogeneous within and between databases. However, there was at least one IR 263 estimated by the algorithms in each database that was comparable with the reference value from the TESSy surveillance system, although some false positives were probably included. 264 265 Based on a few assumptions, that may have overestimated sensitivity, it was estimated that 266 the PPV and sensitivity of the algorithm detecting PC diagnoses in SIDIAP ranged from 54% 267 to 72% and from 83% to 85%, respectively, and that the IRs of Bordetella pertussis in the 268 corresponding population ranged from 25.9 to 35.5 per 100,000 person-years, against the 269 TESSy surveillance system estimate of 21.2.

#### 270 4.1. General comments

Three components were expected to have a high PPV: PC and inpatient specific diagnoses,
and positive laboratory results. Two were expected to have lower PPV (PC and inpatient
unspecified diagnoses). One was expected to be sensitive (prescription of a laboratory test),
two were very unspecific (symptoms and symptoms in infants) and were planned to be used
only in combination with the last component (prescription or use of macrolides) in a 30-days
window of time.

In all the databases, at least one composite algorithm estimated a number of cases that was
compatible with the number expected from the TESSy surveillance system, but this was not
with the combination of the components which was expected to have a high PPV (specific
diagnoses and laboratory tests) in any of the databases. One possible explanation could be that

281 it takes several days to confirm the diagnosis of pertussis after the disease is suspected, and 282 there may be no opportunity for the specific diagnosis to be recorded if the patient does not 283 return to the healthcare facility. Another possible explanation may be that the medical 284 personnel may not see the need to update the record for the purposes of clinical care. This 285 attitude may be influenced by the level of awareness of possible reuse of electronic records 286 for research purposes. These potential explanations may have varying levels of impact in the 287 different databases. For example, in some databases we observed that among the cases 288 detected by a diagnostic component (unspecified or specific), the specific diagnosis was more 289 frequent, indicating that some clinicians might have been more aware about potential research 290 uses of the databases and therefore entered specific diagnoses rather than free text, which was 291 common for unspecified diagnoses.

292 Based on the results of this study, in all the databases it is now possible to design sensitivity 293 analysis using a more specific (but less sensitive) definition of pertussis. In case of 294 heterogeneity in the results of a study on pertussis, designing such sensitivity analyses should 295 be considered as a valid option. On the other hand, in all the databases there is now a possible 296 choice among different sensitive algorithms: we explored several of them, among which 297 'unspecified diagnoses' (the most conservative) and 'test' (the least conservative). Even 298 though these algorithms are likely to have lower PPVs, they may still be useful for sensitivity 299 analyses, especially if there are reasons to think that a specific algorithm could be affected by 300 differential misclassification. For example, pertussis may be more readily suspected and 301 tested for in unvaccinated children, and therefore would be recorded in a more accurate 302 manner.

We developed a component for infants that we though would be sensitive and, although it was likely to have a low PPV, it was less prone to differential misclassification, because it captured symptoms that physicians may not think of as being related with pertussis. However

306 this component proved to be unusable; in reality, when we added a secondary component for 307 concurrent macrolide use there were very few cases that would have been expected to be 308 found in infants with an infection. In contrast, we developed a component specifically for the 309 symptom 'pertussis-like cough' (tos pertusoide in Spanish language) that was apparently 310 specific for pertussis cases that were only found in the BIFAP database. Not only did the 311 majority of cases have a concurrent record of prescription of macrolides, but a manual review 312 of a sample of 100 records including physician free text comments, found 2 cases of 313 unspecified pertussis and 2 cases of suspected pertussis. Therefore, this component may be 314 considered for sensitivity analysis.

## 315 4.2. Compatibility with TESSy and seroprevalence surveys

316 In this study we were able to compare the IRs estimated for paediatric cohorts in five 317 databases using the various algorithms with the national IR estimates from ECDC's TESSy 318 surveillance database. The cases captured by the two types of systems were expected to be 319 slightly different, for various reasons. First, TESSy provides estimates at the national level using census denominators, while three of the databases participating in this study had a 320 321 regional/multiregional scope (SIDIAP, BIFAP and PEDIANET) and two were based on a 322 representative sample of the national population (THIN, RCGP RSC). Therefore it is possible 323 that some clusters of the infectious disease might be under or over-represented in these 324 database. Second, we collected only confirmed cases from TESSy, while some true cases 325 captured by a PC database with a sensitive algorithm may never be confirmed (under 326 ascertainment), or may never be notified (underreporting) [19, 20]. Thus the databases may be 327 a complementary source of true cases which are not notified, while adding potentially false 328 positive cases. Finally, the TESSy data for pertussis may also be affected by under 329 ascertainment and underreporting.

330 The IR found for PEDIANET, which was much higher than the IR estimate from TESSy for 331 Italy (43.9 vs 5.4), may be explained by a combination of both phenomena discussed above. 332 PEDIANET collects data from PC physicians working in the Italian region Veneto, in the 333 North East of the country. The Regional Office for Infectious Diseases of the Veneto Region 334 provided an estimated IR of 10.0 to the data custodians of PEDIANET. This shows that the 335 region had a higher pertussis notification rate than at the national level for 2012 and 2014, although almost all the diagnoses in PEDIANET were unspecified. However, the regional 336 337 estimate could be underestimated because of under ascertainment. Finally, as in the other 338 databases, many cases in PEDIANET could be false positives. In general, if estimates of the 339 PPV of the diagnoses are available, the estimated IR from databases can provide a 340 quantitative estimate of under ascertainment and under notification in TESSy. Vice versa, if 341 under notification to TESSy is known to be small, estimates of the PPV for the algorithm can 342 be obtained. 343 Results from seroprevalence surveys have provided estimates for the incidence of Bordetella

345 Results from seroprevalence surveys have provided estimates for the incidence of *Borderetta* 344 *pertussis* infection [21-23]. These have provided prevalence estimates beyond those of the 345 surveillance systems, partly as they also capture asymptomatic or mildly symptomatic 346 infections. On the contrary, in this study, we observed that estimates of incidence obtained 347 from databases are roughly comparable with those of TESSy.

## 348 4.3. Scope of the component strategy

The scope of this component strategy goes beyond ADVANCE and has the potential of being a comprehensive tool to address heterogeneity and disease misclassification in databases, particularly in multi-database pharmacoepidemiology studies, when the characteristics of the databases affect the operational definition of the outcomes and benchmarking.

353 Inspection of components can provide knowledge that can inform the interpretation of the

354 heterogeneity of the study results. The component strategy can support quantitative bias

355 analysis. In this study, we first developed a set of components with increasing sensitivity and 356 decreasing PPVs. We explored several scenarios for possible PPVs of the components, but in 357 many European databases, estimating directly the PPV of simple algorithms such as 358 components is feasible in a relatively timely and inexpensive way [24-26]. If this is possible, 359 then a consequence of our formulas in Table 3 is that the only value needed to obtain a 360 complete estimate of validity is the sensitivity of the composition of the algorithms, as the rest 361 can be analytically derived. In many cases, sensitivity of the composition could be argued to 362 be very high. In the case of pertussis, we can speculate for instance that cases that were 363 missed from SIDIAP were either seen in a hospital outside of the network that transmits their 364 data to the database, or were very mild and did not require medical attention. The percentage 365 of cases with those characteristics may be estimated from external sources. If estimating this 366 quantity is not possible, then the formulas of Table 3 can still be applied, and they can provide 367 an upper limit for the sensitivity of all the components, that is, the maximum possible 368 sensitivity: to obtain this, it is sufficient to make an assumption on the sensitivity of the 369 *composite* algorithm.

370 If the validity of the variables that enter the analysis can be convincingly proven to be high, 371 this analysis provides evidence that the study results are robust to misclassification. If not, 372 comparing the distribution of components across exposure strata can indicate if differential 373 misclassification is to be suspected. If it is suspected, it can be an important source of bias, as 374 shown by the simulations we report here, as well as in other studies [5, 6]. Components with 375 different validity can, thus, be used to design sensitivity analyses of the study results, applying 376 repeated adjustments for validity to check if the result is robust. If both the PPV and 377 sensitivity are suspected to be non-differential, then the estimate may be unbiased, but the 378 confidence intervals of the estimate need to be adjusted for validity [8]. In future work, the 379 estimates provided by the component strategy could be validated against actual validation

studies. Moreover, the components could be analysed using latent class modelling, whichenables to estimate the validity conditional on various covariates, e.g., age [27].

#### 382 4.4. Strengths and limitations

383 In this study, we used standardised component algorithms as a transparent way of 384 documenting the data extraction process across multiple databases. At the same time, we 385 could also perform a qualitative evaluation of the expected validity of each component 386 Bordetella pertussis, based on its specified semantics and setting. Quantitative scenarios for 387 the validity of each component can also be made using the same approach. We showed that 388 estimates of the validity of various composite algorithms can then be derived in a purely 389 algebraic manner. We could use the incidence estimates based on data from the TESSy 390 surveillance system, which is where European Union member states are required to report 391 pertussis cases, as a reference value, although we cannot exclude the possibility that they may 392 also be subject to under ascertainment and underreporting.

393 The estimates of sensitivity that we obtained for SIDIAP cannot be generalised to the other 394 PC databases. The sensitivity of the PC databases depends on how often a person with the 395 disease symptoms would seek attention in a PC practice. Although in all the databases, the PC 396 physicians have a gatekeeper role, emergency care can be sought without PC referral, and PC 397 practices may not be accessible at night or weekends. Referrals from other settings may be 398 recorded in the PC practice, but no automatic mechanism is in place. In the absence of a 399 database-specific estimate, estimates from another database are a realistic alternative to 400 assuming that sensitivity is 100%.

#### 401 **5. Conclusions**

This study demonstrated the value of quantifying the impact of different event-finding
algorithms across databases and the possibility of benchmarking with disease surveillance
data as well as assessing validity estimates when data from different settings can be linked.

- 405 The validity parameters could be used to correct disease IR estimates from healthcare
- 406 databases.

408	Acknowledgements
409	The authors medical writing and editorial assistance from Margaret Haugh (MediCom
410	Consult, Villeurbanne, France).
411	
412	

## 413 **Disclaimer**

414 The results described in this publication are from the proof of concept studies conducted as 415 part of the IMI ADVANCE project with the aim of testing the methodological aspects of the 416 design, conduct and reporting of studies for vaccine benefit-risk monitoring activities. The 417 results presented relate solely to the methodological testing and are not intended to inform 418 regulatory or clinical decisions on the benefits and risks of the exposures under investigation. 419 This warning should accompany any use of the results from these studies and they should be 420 used accordingly. The views expressed in this article are the personal views of the authors and 421 should not be understood or quoted as being made on behalf of or reflecting the position of 422 the agencies or organisations with which the authors are affiliated. 423

## 425 **Funding source**

- 426 The Innovative Medicines Initiative Joint Undertaking funded this project under ADVANCE
- 427 grant agreement n° 115557, resources of which were composed of a financial contribution
- 428 from the European Union's Seventh Framework Programme (FP7/2007-2013) and in kind
- 429 contributions from EFPIA member companies.

430

## **Declaration of potential conflicts of interest**

Rosa Gini, Caitlin Dodd, Kaatje Bollaerts, Claudia Bartolini, Giuseppe Roberto, Consuelo Huerta-Alvarez, Elisa Martín-Merino, Talita Duarte-Salles, Gino Picelli, Lara Tramontan, Giorgia Danieli, Ana Correa, Chris McGee, Benedikt Becker, Charlotte Switzer, Jorgen Bauwens, Nicoline van der Maas, Gianfranco Spiteri, Emmanouela Sdona declared no conflicts of interest. Sonja Gandhi-Banga declared that she works for Sanofi Pasteur and holds company shares. Daniel Weibel declared that he has received personal fees from GSK for work unrelated to the submitted work. Miriam Sturkenboom declared that she has received grants from Novartis, CDC and Bill & Melinda Gates Foundation for work unrelated to the submitted work.

### 445 **References**

- 446 [1] Sturkenboom M, van der Aa L, Bollaerts K, Emborg HD, Ferreira G, Gino R, et al. The
- 447 ADVANCE distributed network system for evidence generation on vaccines coverage,
- 448 benefits and risks based on electronic health care data. Vaccine. 2018;Paper 2 in supplement.
- [2] Sturkenboom M, Weibel D, van der Aa L, Braeye T, Gheorge M, Becker B, et al.
- 450 ADVANCE database characterization and fit for purpose assessment for multi-country studies
- 451 on the coverage, benefits and risks of vaccinations. Vaccine. 2018;Paper 3 in Supplement.
- 452 [3] Gini R, Schuemie M, Brown J, Ryan P, Vacchi E, Coppola M, et al. Data extraction and
- 453 management in networks of observational health care databases for scientific research: a
- 454 comparison of EU-ADR, OMOP, Mini-Sentinel and MATRICE strategies. EGEMS
- 455 (Washington, DC). 2016;4:1189.
- 456 [4] Avillach P, Coloma PM, Gini R, Schuemie M, Mougin F, Dufour JC, et al. Harmonization
- 457 process for the identification of medical events in eight European healthcare databases: the
- 458 experience from the EU-ADR project. J Am Med Inform Assoc. 2013;20:184-92.
- 459 [5] Funk MJ, Landi SN. Misclassification in administrative claims data: quantifying the
- 460 impact on treatment effect estimates. Curr Epidemiol Rep. 2014;1:175-85.
- 461 [6] Hofler M. The effect of misclassification on the estimation of association: a review. Int J
- 462 Methods Psychiatr Res. 2005;14:92-101.
- 463 [7] De Smedt T, Merrall E, Macina D, Perez-Vilar S, Andrews N, Bollaerts K. Bias due to
- 464 differential and non-differential disease- and exposure misclassification in studies of vaccine
- 465 effectiveness. PloS One. 2018;13:e0199180.
- 466 [8] Brenner H, Gefeller O. Use of the positive predictive value to correct for disease
- 467 misclassification in epidemiologic studies. Am J Epidemiol. 1993;138:1007-15.

- 468 [9] Roberto G, Leal I, Sattar N, Loomis AK, Avillach P, Egger P, et al. Identifying cases of
- 469 type 2 diabetes in heterogeneous data sources: strategy from the EMIF Project. PloS One.
  470 2016;11:e0160648.
- 471 [10] Barlow RS, Reynolds LE, Cieslak PR, Sullivan AD. Vaccinated children and adolescents
- 472 with pertussis infections experience reduced illness severity and duration, Oregon, 2010-2012.
- 473 Clin Infect Dis. 2014;58:1523-9.
- 474 [11] McNamara LA, Skoff T, Faulkner A, Miller L, Kudish K, Kenyon C, et al. Reduced
- 475 severity of pertussis in persons with age-appropriate pertussis vaccination-United States,
- 476 2010-2012. Clin Infect Dis. 2017;65:811-8.
- 477 [12] European Parliament and of the Council. (Decision EU 2012) Commission implementing
- 478 decision of 8 August 2012 amending Decision 2002/253/EC laying down case definitions for
- 479 reporting communicable diseases to the Community network under Decision No 2119/98/EC
- 480 of the European Parliament and of the Council. Annex to L 262. Official Journal of the
- 481 European Union 27/9/2012. Available at: <u>http://eur-lex.europa.eu/legal-</u>
- 482 content/EN/TXT/PDF/?uri=CELEX:32012D0506&qid=1428573336660&from=EN#page=22
- 483 . Accessed on: 9 November 2018.
- 484 [13] European Centre for Disease Prevention and Control. The European Surveillance System
- 485 (TESSy). Available at: <u>https://ecdc.europa.eu/en/publications-data/european-surveillance-</u>
- 486 <u>system-tessy</u>. Accessed on: 9 November 2018.
- 487 [14] Becker BFH, Avillach P, Romio S, van Mulligen EM, Weibel D, Sturkenboom M, et al.
- 488 CodeMapper: semiautomatic coding of case definitions. A contribution from the ADVANCE
- 489 project. Pharmacoepidemiol Drug Saf. 2017;26(8):998-1005.
- 490 [15] Bellettini CV, de Oliveira AW, Tusset C, Baethgen LF, Amantea SL, Motta F, et al.
- 491 [Clinical, laboratorial and radiographic predictors of Bordetella pertussis infection]. Revista
- 492 paulista de pediatria : orgao oficial da Sociedade de Pediatria de Sao Paulo. 2014;32:292-8.

- 493 [16] Hurtado-Mingo A, Mayoral-Cortes JM, Falcon-Neyra D, Merino-Diaz L, Sanchez-
- 494 Aguera M, Obando I. [Clinical and epidemiological features of pertussis among hospitalized
- 495 infants in Seville during 2007-2011]. Enferm Infecc Microbiol Clin. 2013;31:437-41.
- 496 [17] Eurostat. Population data. Available at: <u>https://ec.europa.eu/eurostat/web/population-</u>
- 497 <u>demography-migration-projections/population-data/database</u>. Accessed on: 12 November
- 498 2018.
- [18] Ulm K. A simple method to calculate the confidence interval of a standardized mortality
  ratio (SMR). Am J Epidemiol. 1990;131:373-5.
- 501 [19] McDonald SA, Teunis P, van der Maas N, de Greeff S, de Melker H, Kretzschmar ME.
- 502 An evidence synthesis approach to estimating the incidence of symptomatic pertussis
- 503 infection in the Netherlands, 2005-2011. BMC Infect Dis. 2015;15:588.
- 504 [20] Schielke A, Takla A, von Kries R, Wichmann O, Hellenbrand W. Marked underreporting
- 505 of pertussis requiring hospitalization in infants as estimated by capture-recapture
- 506 methodology, Germany, 2013-2015. Pediatr Infect Dis J. 2018;37:119-25.
- 507 [21] Barkoff AM, Grondahl-Yli-Hannuksela K, He Q. Seroprevalence studies of pertussis:
- 508 what have we learned from different immunized populations. Pathog Dis. 2015;73.
- 509 [22] de Greeff SC, de Melker HE, van Gageldonk PG, Schellekens JF, van der Klis FR,
- 510 Mollema L, et al. Seroprevalence of pertussis in The Netherlands: evidence for increased
- 511 circulation of Bordetella pertussis. PloS One. 2010;5:e14183.
- 512 [23] de Melker HE, Versteegh FG, Schellekens JF, Teunis PF, Kretzschmar M. The incidence
- 513 of Bordetella pertussis infections estimated in the population from a combination of
- 514 serological surveys. J Infect. 2006;53:106-13.
- 515 [24] Coloma PM, Valkhoff VE, Mazzaglia G, Nielsson MS, Pedersen L, Molokhia M, et al.
- 516 Identification of acute myocardial infarction from electronic healthcare records using different
- 517 disease coding systems: a validation study in three European countries. BMJ Open. 2013;3.

- 518 [25] Valkhoff VE, Coloma PM, Masclee GM, Gini R, Innocenti F, Lapi F, et al. Validation
- 519 study in four health-care databases: upper gastrointestinal bleeding misclassification affects
- 520 precision but not magnitude of drug-related upper gastrointestinal bleeding risk. J Clin
- 521 Epidemiol. 2014;67:921-31.
- 522 [26] Gini R, Schuemie MJ, Mazzaglia G, Lapi F, Francesconi P, Pasqua A, et al. Automatic
- 523 identification of type 2 diabetes, hypertension, ischaemic heart disease, heart failure and their
- 524 levels of severity from Italian General Practitioners' electronic medical records: a validation
- 525 study. BMJ Open. 2016;6:e012413.
- 526 [27] Hui SL, Walter SD. Estimating the error rates of diagnostic tests. Biometrics.
- 527 1980;36:167-71.

## 529 Figure caption

## 530 Figure 1: Study results for the incidence of tested component and composite algorithms.

531 For each component algorithm the incidence rate per 100,000 person-years is shown. For the

532 composite algorithms, the incidence rates were stratified per type of case: cases detected only

- 533 by the left-hand component (indicated in the label before the Boolean operator 'OR'), cases
- 534 detected by both components, and cases detected by the right-hand component (indicated in
- 535 the label after the key Boolean operator word 'OR'). The dashed line represents the national
- 536 incidence rate per 100,000 person-years based on data from the TESSy surveillance system.
- 537 Data for years 2012 and 2014 were pooled.

## 539 Appendix: Members of ADVANCE consortium (October 2018)

- 540 Full partners
- 541 AEMPS: Agencia Española de Medicamentos y Productos Sanitarios (www.aemps.es)
- 542 ARS-Toscana: Agenzia regionale di sanità della Toscana (https://www.ars.toscana.it/it/)
- 543 ASLCR: Azienda Sanitaria Locale della Provincia di Cremona (www.aslcremona.it)
- 544 AUH: Aarhus Universitetshospital (kea.au.dk/en/home)
- 545 ECDC: European Centre of Disease Prevention and Control (www.ecdc.europa.eu)
- 546 EMA: European Medicines Agency (www.ema.europa.eu)
- 547 EMC: Erasmus Universitair Medisch Centrum Rotterdam (www.erasmusmc.nl)
- 548 GSK: GlaxoSmithKline Biologicals (www.gsk.com)
- 549 IDIAP: Jordi Gol Fundació Institut Universitari per a la Recerca a l'Atenció Primària de Salut
- 550 Jordi Gol i Gurina (http://www.idiapjordigol.com)
- 551 JANSSEN: Janssen Vaccines Prevention B.V. (http://www.janssen.com/infectious-diseases-
- and-vaccines/crucell)
- 553 KI: Karolinska Institutet (ki.se/meb)
- 554 LSHTM: London School of Hygiene & Tropical Medicine (www.lshtm.ac.uk)
- 555 MHRA: Medicines and Healthcare products Regulatory Agency (www.mhra.gov.uk/)
- 556 MSD: Merck Sharp & Dohme Corp. (www.merck.com)
- 557 NOVARTIS: Novartis Pharma AG (www.novartisvaccines.com)
- 558 OU: The Open University (www.open.ac.uk)
- 559 P95: P95 (www.p-95.com)
- 560 PEDIANET: Società Servizi Telematici SRL (www.pedianet.it)
- 561 PFIZER: Pfizer Limited (www.pfizer.co.uk)
- 562 RCGP: Royal College of General Practitioners (www.rcgp.org.uk)
- 563 RIVM: Rijksinstituut voor Volksgezondheid en Milieu (www.rivm.nl)

- 564 SCIENSANO: Sciensano (https://www.sciensano.be)
- 565 SP: Sanofi Pasteur (www.sanofipasteur.com)
- 566 SSI: Statens Serum Institut (www.ssi.dk)
- 567 SURREY: The University of Surrey (www.surrey.ac.uk)
- 568 SYNAPSE: Synapse Research Management Partners, S.L. (www.synapse-managers.com)
- 569 TAKEDA: Takeda Pharmaceuticals International GmbH (www.tpi.takeda.com)
- 570 UNIBAS-UKBB: Universitaet Basel Children's Hospital Basel (www.unibas.ch)
- 571 UTA: Tampereen Yliopisto (www.uta.fi)

## 572 Associate partners

- 573 AIFA: Italian Medicines Agency (www.agenziafarmaco.it)
- 574 ANSM: French National Agency for Medicines and Health Products Safety (ansm.sante.fr)
- 575 BCF: Brighton Collaboration Foundation (brightoncollaboration.org)
- 576 EOF: Helenic Medicines Agency, National Organisation for Medicines (www.eof.gr)
- 577 FISABIO: Foundation for the Promotion of Health and Biomedical Research
- 578 (www.fisabio.es)
- 579 HCDCP: Hellenic Centre for Disease Control and Prevention (www.keelpno.gr)
- 580 ICL: Imperial College London (www.imperial.ac.uk)
- 581 IMB/HPRA: Irish Medicines Board (www.hpra.ie)
- 582 IRD: Institut de Recherche et Développement (www.ird.fr)
- 583 NCE: National Center for Epidemiology (www.oek.hu)
- 584 NSPH: Hellenic National School of Public Health (www.nsph.gr)
- 585 PHE: Public Health England (www.gov.uk/government/organisations/public-health-england)
- 586 THL: National Institute for Health and Welfare (www.thl.fi)
- 587 UMCU: Universitair Medisch Centrum Utrecht (www.umcu.nl)
- 588 UOA: University of Athens (www.uoa.gr)

- 589 UNIME: University of Messina (www.unime.it)
- 590 Vaccine.Grid: Vaccine.Grid (http://www.vaccinegrid.org/)
- 591 VVKT: State Medicines Control Agency (www.vvkt.lt)
- 592 WUM: Polish Medicines Agency Warszawski Uniwersytet Medyczny
- 593 (https://wld.wum.edu.pl/)

## Table 1

# Table 1: Sets of concepts selected for the component algorithms

Each set of concepts can contain one or more concepts, each described and, if available, with a Concept Unique Identifier of the Unified Medical

Language System.

			Concept
Concept set label	Concept set description	Concept	Unique
			Identifier
		Bordetella pertussis	C0043167
		Whooping cough due to	
(Bordetella pertussis)		Bordetella pertussis without	C2887068
		pneumonia	
		Whooping cough due to	
	Concepts referring to diagnoses specifically mentioning pertussis	Bordetella pertussis with	C2887069
	induced by an infection of Bordetella pertussis	pneumonia	
	Concepts referring to diagnoses which refer to pertussis, but without	a Whooping cough due to	C0042169
(Pertussis unspecified)	specific indication that Bordetella pertussis is responsible for the	unspecified organism	C0043108
	infection	Bordetella infections	C0006015

Concept set label	Concept set description	Concept	Concept			
		Whooping cough-like syndrome	C0343485			
		Notification of whooping cough				
		Pneumonia in pertussis	C0155865			
(9	This set of concepts was introduced because the Spanish translation of	of				
with pertussis)	'whooping cough' was found to be considered by physicians as a	Concept of 'tos pertusoide' in				
	symptom, not as a diagnosis	Spanish general practice				
		Apnea	C0003578			
(Symptoms in infants)		Cyanosis	C0010520			
	Concepts referring to symptoms that were found to be predictive of	Post-tussive vomiting	C1740793			
	pertussis in infants [13, 14]	Paroxysms of coughing	C0231911			
(Macrolides)	Use of macrolides	Macrolides				
		Polymerase chain reaction test				
	The concepts listed in this set indicate the prescription of tests that an	reCulture or serology				
(Bordetella pertussis test	considered to be confirmatory of a <i>Bordetella pertussis</i> infection	Isolation of Bordetella pertussis				
		from a clinical specimen				

Concept set label	Concept set description	Concept	Concept					
	Positive							
		reaction test						
(Positive result from a	The concepts listed in this set indicate a positive result from a tests	Positive culture or serology						
Bordetella pertussis test)	confirmatory of a Bordetella pertussis infection	Positive isolation of Bordetella						
	pertussis from a clinical							
		specimen						

# Table 2: Components for pertussis

The concept sets referred to by the words in round parentheses can be found in Table 1.

Name	Setting	Data domain	Concept set
PC diagnosis, specific	Primary care practice	Diagnosis	(Bordetella pertussis)
Inpatient diagnosis, specific	Hospital	Diagnosis	(Bordetella pertussis)
PC diagnosis, unspecified	Primary care practice	Diagnosis	(Pertussis unspecified)
Inpatient diagnosis, unspecified	Hospital	Diagnosis	(Pertussis unspecified)
Symptoms	Primary care practice	Diagnosis or signs/symptoms	(Symptoms compatible with pertussis)
Symptoms in infants	Primary care practice	Diagnosis or signs/symptoms	(Symptoms in infants)

Name	Setting	Data domain	Concept set						
Test	Any setting where a test can be prescribed, or facility where the test is administered	Laboratory test	(Bordetella pertussis test)						
Positive laboratory results	Any setting where a health professional records the results of a test, or facility where the results of the test are generated	Results from laboratory test	(Positive result from a <i>Bordetella pertussis</i> test)						
Drug use	Facility dispensing medications or primary care practice issuing prescriptions	(Macrolides)							
Secondary components									
Symptoms and drugs within 30 days	Symptoms and drugsA patient is positive if they have both a record of symptoms and of drug use, and the interval between the dates iswithin 30 daysless than 30 days								
Symptoms in infants and drugs within 30 days	A patient is positive if they are 0 or 1 and has both a record of symptoms in infants and of drug use, and the interval between the dates is less than 30 days								

**Table 3:** Analytic formulae linking the true proportion of pertussis and validity indices of one or two algorithms

In the formulas,  $\Pi$  is the true proportion of cases of pertussis, P is the proportion of cases detected by the algorithm, SE is the sensitivity and PPV is the positive predictive value of the algorithm.

Known parameters	Formula to derive another parameter
One algorithm	
PPV and SE	$\pi = \frac{P \times PPV}{SE}$
PPV and Π	$SE = \frac{P \times PPV}{\pi}$
SE and Π	$PPV = \frac{SE \times \pi}{P}$
Two algorithms A and B	
SE of A, of B, and of A	$SE_A \operatorname{or}_B = SE_A + SE_B - SE_A \operatorname{and}_B$
AND B	
$\Pi$ and PPV of A, of B, and	$SE_{A \text{ or } B} = \frac{P_A \times PPV_A}{\pi} + \frac{P_B \times PPV_B}{\pi} - \frac{P_A \text{ and } B \times PPV_A \text{ and } B}{\pi}$
of A AND B	
SE of A <b>OR</b> B, and PPV of	$\pi = \frac{P_A \times PPV_A + P_B \times PPV_B - P_A \text{ and } B \times PPV_A \text{ and } B}{2P}$
A, of B, and of A <b>AND</b> B	$SE_A$ or $_B$
PPV of A, of B, and of A	$PPV_{A \text{ or } B} = \frac{P_A \times PPV_A + P_B \times PPV_B - P_{A \text{ and } B} \times PPV_A \text{ and } B}{P_A \text{ or } B}$
AND B	

**Table 4:** Study results. Number of person-years (PYs) entering the study in each database. Incidence rates of pertussis per 100,000 children aged 0-14, with 95% confidence interval (CI), from the TESSy surveillance system in the corresponding country and the estimate incidence rate per 100,000 for each component algorithm are shown. In composite algorithms, the incidence rates were stratified per type of case: cases detected only by the left-hand component (indicated in the label before the keyword 'OR'), cases detected by both components, and cases detected by the right-hand component (indicated in the label after the keyword 'OR'). Data for years 2012 and 2014 were pooled.

		1	/			U							/							
DB	SIDIAP (Spain)		BIFAP (Spain)			PEDIANET (Italy)			THIN (United Kingdom)				RCGP (United Kingdom)							
Person-years	488,847			796,324			88, 754			1,387,939				411,404						
TESSy (IR and	21.2 (20.5-22.0)				21.2 (20.5-22.0)				5.4 (5	.1-5.8)			13.4 (1)	3.0-13.9)			13.4 (13	3.0-13.9)		
95% ČÌ)																				
								Component	algorithms (N	and IR per 1	00,000 PYs)									
PC diagnosis,	21 (4.3)				0 (0.0)				2 (2.3)				79 (5.7)				21 (5.1)			
specific																				
PC diagnosis,	173 (35.4)				135 (17.0)				37 (41.7)				178				77 (18.7)			
unspecified													(12.8)							
Inpatient	27 (5.5)																			
diagnosis,																				
specific								_								-				
Inpatient	26 (5.3)																			
diagnosis,																				
Grand					1(((20.0)														-	
Symptoms					100 (20.8)														-	
drug within					122 (15.3)															
30days																				
Symptoms in	27 (5 5)	1					1		6(68)				172			1	30 (7 3)		1	
infants	27 (0.0)								0 (0.0)				(12.4)				00(1.0)			
Symptoms in	1 (0.2)												8 (0.6)							
infants and																				
drug within																				
30days																				
Test	96 (19.6)				38 (4.8)				38 (42.8)				209				32 (7.8)			
D	10 (2.0)				0 (0 0)								(15.1)							
Positive	19 (3.9)				0 (0.0)								3 (0.2)							
regulte																				
Tesuits						1			Composito	algorithms										
	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)	N (IP)
	N (IN)	in left-	In hoth	In right-hand	iv (iiv)	in left-	In hoth	In right-	N (IN)	in left-	In hoth	In right-	in (inc)	in left-	In hoth	In right-	in (inc)	in left-	In hoth	In right-
		hand	compo	component		hand	compo	hand		hand	compo	hand		hand	compo	hand		hand	compo	hand
		componen	nets	only		compone	nets	compone		compone	nets	compone		compone	nets	compone		compon	nets	compone
		t only		- 5		nt only		nt only		nt only		nt only		nt only		nt only		ent only		nt only
PC specific OR	194 (39.6)	21 (4.3)	0 (0.0)	173 (35.4)	135 (17.0)	0 (0.0)	0 (0.0)	135	39 (43.9)	2 (2.2)	0 (0.0)	37 (41.7)	246	68 (4.9)	11 (0.8)	167	91 (22.1)	14 (3.4)	7 (1.7)	70 (17.0)
unspecified								(17.0)					(17.7)			(12.0)				
diagnosis																				
Inpatient	52 (10.6)	25 (5.1)	1 (0.2)	26 (5.3)																
specific OR																				
unspecified																				
diagnosis	220 (45 0)	1(0(24.2)	26 (5.2)	26 (5.2)			-													
PC OK Inpatient	220 (45.0)	168 (34.3)	26 (5.3)	26 (5.3)																
DC diagnosis OP	271 (55.4)	77 (15.9)	10 (2.0)	175 (25.9)	160 (21.1)	22 (4 1)	F (0.6)	120	69 (777)	20 (22 0)	8 (0,0)	21 (24.0)	126	101	20 (2 1)	217	115 (29.0)	24 (5.9)	9(10)	92 (20 2)
test	271 (33.4)	//(15.0)	19 (3.9)	175 (55.0)	100 (21.1)	33 (4.1)	5 (0.0)	(163)	09(77.7)	30 (33.8)	0 (9.0)	31 (34.9)	(30.7)	(13.0)	27 (2.1)	(15.6)	115 (20.0)	24 (5.0)	0 (1.9)	03 (20.2)
Positive lab	197 (40 3)	3(0.6)	16 (3 3)	178 (36.4)	135 (17.0)	0 (0 0)	0 (0 0)	135					247	1 (0 1)	2 (0 1)	244				
results OR PC	. ()	. ()	. (a.a.)			, (0.0)	. ()	(17.0)					(17.8)	- ()	. ()	(17.6)				
diagnosis							1						,		1	,				
PC diagnosis OR					255 (32.0)	133	2 (0.3)	120							1					
symptoms and						(16.7)		(15.1)							1	1				
drugs																				
Any diagnosis	223 (45.6)	204 (41.7)	16 (3.3)	3 (0.6)																
OR positive lab							1	1							1	1				
results						1	1	1						1	1	1				

Table 4



Figure 1

Supplemental File 1 Click here to download Supplemental Files: Supplementary File 1.docx Supplemental File 2 Click here to download Supplemental Files: Supplementary\_File\_2.pdf Supplemental Table 1 Click here to download Supplemental Files: Supplementary Table 1.docx Supplemental Table 2 Click here to download Supplemental Files: Supplementary Table 2.pdf Supplemental Figure 1 Click here to download Supplemental Files: Supplementary Figure 1.docx