

Research article

An instrument to identify computerised primary care research networks, genetic and disease registries prepared to conduct linked research: TRANSFoRm International Research Readiness (TIRRE) survey

Cite this article: Jennings E, de Lusignan S, Michalakidis G, Krause P, Sullivan F, Liyanage H, Delaney BC. An instrument to identify computerised primary care research networks, genetic and disease registries prepared to conduct linked research: TRANSFoRm International Research Readiness (TIRRE) survey. *J Innov Health Inform.* 2018;25(4):207–220.

<http://dx.doi.org/10.14236/jhi.v25i4.964>

Copyright © 2018 The Author(s). Published by BCS, The Chartered Institute for IT under Creative Commons license <http://creativecommons.org/licenses/by/4.0/>

Author address for correspondence:

Simon de Lusignan
Professor
Primary Care and Clinical Informatics
Department of Clinical and Experimental Medicine
University of Surrey
Guildford GU2 7XH, UK
Email: s.lusignan@surrey.ac.uk

Accepted December 2018

Emily Jennings

Hon. Research Assistant, Department of Clinical and Experimental Medicine, University of Surrey, Guildford, UK

Simon de Lusignan

Professor, Primary Care and Clinical Informatics, Department of Clinical and Experimental Medicine, University of Surrey, Guildford, UK

Georgios Michalakidis

PhD Data Analytics, Department of Computer Science, University of Surrey, Guildford, UK

Paul Krause

Professor, Complex Systems, Department of Computer Science, University of Surrey, Guildford, UK

Frank Sullivan

Professor, Primary Care, Population Health Sciences, University of Dundee, Dundee, UK

Harshana Liyanage

Research Fellow, Department of Clinical and Experimental Medicine, University of Surrey, Guildford, UK

Brendan C. Delaney

Professor, Primary Care Research, Department of Surgery and Cancer, Imperial College London, St Mary's Campus, London, UK

ABSTRACT

Purpose The Translational Research and Patients safety in Europe (TRANSFoRm) project aims to integrate primary care with clinical research whilst improving patient safety. The TRANSFoRm International Research Readiness survey (TIRRE) aims to demonstrate data use through two linked data studies and by identifying clinical data repositories and genetic databases or disease registries prepared to participate in linked research.

Method The TIRRE survey collects data at micro-, meso- and macro-levels of granularity; to fulfil data, study specific, business, geographical and readiness requirements of potential data providers for the TRANSFoRm demonstration studies. We used descriptive statistics to differentiate between demonstration-study compliant and non-compliant repositories. We only included surveys with >70% of questions answered in our final analysis, reporting the odds ratio (OR) of positive responses associated with a demonstration-study compliant data provider.

Results We contacted 531 organisations within the European Union (EU). Two declined to supply information; 56 made a valid response and a further 26 made a partial response. Of the 56 valid responses, 29 were databases of primary care data, 12 were genetic databases and 15 were cancer registries. The demonstration compliant primary care sites made 2098 positive responses compared with 268 in non-use-case compliant data sources [OR: 4.59, 95% confidence interval (CI): 3.93–5.35, $p < 0.008$]; for genetic databases: 380:44 (OR: 6.13, 95% CI: 4.25–8.85, $p < 0.008$) and cancer registries: 553:44 (OR: 5.87, 95% CI: 4.13–8.34, $p < 0.008$).

Conclusions TIRRE comprehensively assesses the preparedness of data repositories to participate in specific research projects. Multiple contacts about hypothetical participation in research identified few potential sites.

Keywords: medical informatics, family practice, medical records systems, electronic health records, diabetes mellitus, Barrett's disease

ABBREVIATIONS

IT - information technology; EMR - electronic medical records; TRANSFoRm - Translational Research and Patients safety in Europe; TIRRE survey - TRANSFoRm International Research Readiness survey; IBM SPSS - IBM Statistical Package for Social Sciences; eHR - electronic health record; OR - odds ratio; ICD - International Classification of Disease; ICPC - International Classification of Primary Care; SNOMED - Systematised Nomenclature of Medicine; CTv3 - Clinical Terms Version 3; ATC - Anatomical Therapeutic Chemical; HL7 - Health Level-7; RIM - Reference Information Model; CDISC - Clinical Data Interchange Standards Consortium; BRIDG - Biomedical Research Integrated Domain Group; CSV - Comma Separated Values; CPRD - Clinical Practice Research Datalink.

INTRODUCTION

Large databases of health data are widely used for research but less often combined.¹ Linked data facilitates better measurement of clinical performance and patient health outcomes in health care systems.² Technical challenges of linking data are mostly considered to be the key barrier of integrating disparate heterogeneous data sources.³ Data privacy legislations can considerably hinder research in a multinational setting.⁴ Data collected within primary care have been computerised since the 1990s⁵ with data widely used for research,⁶ but with relatively little linkage of data beyond disease-specific programmes in individual localities. In the United States, the federal electronic medical records mandate aims not only to save money but also to modernise health information technology (IT). A team of RAND Corporation researchers projected in 2005 that a move towards health IT could potentially save \$81 billion. However, this saving has far from materialised and despite the recommendations, spending in the US has increased over the past 9 years by \$800 billion.⁷ The increase in spending was, in part, attributed to the slow adoption of health IT systems that are neither interoperable nor easy to use.

The Translational Research and Patient Safety in Europe (TRANSFoRm) project aims to reduce barriers to conducting research using routine healthcare data across Europe.^{8–10} The European eHealth Action Plan prioritises interoperability between health records so that internationally comparable data can be collected on the quality of care and for research.¹¹ The TRANSFoRm International Research Readiness (TIRRE) survey was developed and designed to collect

information about these data sources with the primary aim of assessing the preparedness of disease registries, throughout Europe, to conduct linked research using the TRANSFoRm project (Appendix 1). The TRANSFoRm requirements for the TIRRE instrument were that it could assess the feasibility of conducting two simulated studies (use-cases): one on the genetics of response to oral anti-diabetic medication; the other on the relationship between anti-indigestion medication, Barrett's disease, oesophageal cancer and the quality of life. The 'use-cases' were designed to capture how primary care recorded oesophageal reflux might be a prodrome of cancer; and any genetic predisposition to complications of people with type 2 diabetes.⁶

Initially, we focussed on whether a repository had the required dataset 'data readiness' or the technical mechanism to extract data in the format needed for analysis 'record readiness'. The different levels are ascertained by reviewing the metadata captured from the survey. The overall assessment gives an indication as to whether a database is capable of contributing to linked research (*i.e.* 'linked research readiness'). However, a stakeholder analysis, expert opinion¹² and pilot data collection,¹³ all pointed toward the necessity to take a broader approach and to include socio-cultural and business process aspects of readiness. We consequently developed a more systematic approach to analysing the requirements for generic data linkage studies, as well as those specific to the TRANSFoRm use-cases.¹⁴ Our pilot TIRRE survey indicated that this instrument identified databases that could be used to conduct the more sophisticated translational and biomedical research planned within the TRANSFoRm project. Consequently, the final version of TIRRE collects data at three levels of granularity – the micro or data level; the meso- or record system level and the macro or health service level. TIRRE also includes study specific questions defined from the oesophagitis and diabetes use-cases. This study reports the results of the first TIRRE survey.

METHOD

Sampling and data collection

Our initial contact was to the health ministry of each EU country and to National Primary Care Organisations. Subsequent strategies included trying to identify sites through Internet and Medline searches, and snowball sampling through contacts made or work references. We also contacted National and European informatics and research networks. We identified sites across Europe willing to participate in the survey by contacting them through email or web-form and we then followed this up with a phone call. We exported these data from the

completed online questionnaires directly into either Microsoft Excel or into Statistical Package for Social Sciences (IBM SPSS). We categorised 'non-compliance' as a respondent who partially completed the online survey, answering <70% of the questions; or as a respondent with whom we had made telephone contact initially was unavailable for their telephone interview or failed to proceed to online completion of the survey. A major component of the workload in this project involved identifying potential survey respondents.

Micro-, meso- and macro-level

The broad scope of the survey emerged from a series of workshops and is composed of a wide range of questions designed to assess how data might be linked, the data itself, extraction methods and social and organisational influences.^{15,16} The final instrument contained 160 questions divided into a framework which consisted of micro-, meso-, macro- and study-specific levels.

- The first section covered micro-level issues and was concerned with the data source, the data itself, metadata, the potential for linkage or achieving semantic interoperability between data sources¹⁷ and details of how many studies have been published using the data.
- The meso-level explored the data extraction,¹⁸ the architecture for the computerised medical record and other data repositories,¹⁹ audit trails and the size of the database.
- The macro-issues related to the nature of the health system, socio-cultural factors and issues relating to the funding, purpose and restrictions on the use of the data.
- Study-specific questions make up the final part of the survey instrument (Supplementary data file, Table S1), these were designed to identify sites that were eligible to participate in the use-cases in pairs of primary care and genetic, or primary care and cancer registry data.

We described the coding systems used to store data, including drug dictionaries and any standards used (the aim was to determine whether there were a small number of possible combinations of coded data to identify within data repositories and the mechanisms for achieving interoperability), the number and details of eHR vendors, vendors of communications and data processing applications routinely used (including their international scope, coding systems offered and if they had common data export formats) along with organisational, policy, cultural or legislative restrictions on data reuse.

Use-case specific

We analysed the process of conducting two use-cases and defined the studies using a framework which defined the micro-, meso- and macro-levels of data and process information required to conduct successful linked research, where multiple data sources are semantically integrated. We

summarised the sites eligible to participate in the use-cases in pairs of primary care and genetic or primary care and cancer registry data. If the database can support a use-case, we consider the site as a use-case compliant site and if it cannot support, we define the site as a non-use-case site. Registries were only eligible if they provided a valid response to the questionnaire. We required as much of the survey to be completed as possible as each part of it was determined from our requirements analysis. We defined a valid response to be one which answered >70% of the questions. Key compulsory answer questions which defined compliance provided information such as valid contact details, a link to another dataset, size of the dataset, data model and details of the coding system, the likely lead time in any approval process and that they have use-case variables available. All sections of the questionnaire provided significant and useful information to determine if the database was use-case specific.

Reporting and analysis

We compared the responses from databases that proved eligible to participate in the use-cases with those who were not. We wanted to explore whether it was more likely that those associated with eligibility would give a positive response to questions than those who were not deemed eligible. A valid response provided by the respondent is considered a positive response. The purpose of this exercise was to identify any questions that were not purposeful and to reduce the number of questions. We identified and reviewed any questions that were not answered positively by any of the use-case eligible respondents on the basis that they were not discriminatory of eligibility to participate in either of the studies.

Statistical methods

We used descriptive statistics (*i.e.* measures of frequency) to describe response rates and quote odds-ratios (ORs), 95% confidence intervals (CIs) and used tests of proportion to report whether sections of the questionnaire helped to discriminate between those able to conduct the use-case or not.

Ethics statement

There was no formal ethics board review. This survey only seeks to report information about the capacity and capability of information sources to be combined to conduct research studies and does not involve any access to personal data. However, the TIRRE survey does check whether data sources collect individual consent and if they contain strong identifiers and if there are restrictions on the use of data.

RESULTS

Sample and data collection for use-case specific defined studies

We made many contacts but received few responses. We contacted 531 different organisations, and later individuals in EU countries (including eHR vendors) and received 56 valid responses. Of the health ministries we contacted, seven provided useful information and a further five responded but could

not provide any helpful information. Only two site representatives declined to participate at this stage (Supplementary data file, Table S2).

eHR vendors

We also collected details of the national or international eHR vendors with a significant presence in one or more EU countries. We contacted 17 companies identified initially, as well as any reported by survey respondents. Nine of these eHR vendors had a presence in more than one country. Two of those contacted started to complete the TIRRE survey instrument but failed to complete the questionnaire. We also approached nine vendors listed by site representatives who completed the questionnaire but they once again expressed no interest in participating in the survey. They did suggest they might consider completing the survey in the future if and when we had something more definite to offer. Few vendors responded; however, when they did reply to the survey, their responses to the questions posed provided useful detail.

Telephone and online completion of the survey

Of the 531 organisations we made contact with, 45 respondents commenced but did not complete the TIRRE survey online (Supplementary data file, Table S3) and 26 made a partial response during telephone enquiries but were then either unavailable for their telephone interview or failed to go ahead and complete the online survey. The initial telephone interviews took 1.5 hours and with experience still took 50–75 minutes. The feedback from the pilot survey suggested that the process took too long and that there was very little incentive for the respondent for completing the survey. While this drawback of the survey could have possibly caused a bias for the responses collected, we consider this as a valuable learning to consider in similar database profiling activities conducted in the future.

Completion of the survey

The valid surveys were on an average returned with 76% of the questions completed and this was consistent across the three respondent groups. Looking at the survey by category, the *Data source* and *Record system* sections were the only ones that fell below the 75% level (many sections were returned with above 90% of the questions completed). The main reason for this was the variation in the skip logic for individual respondents in these sections of the questionnaire (Supplementary data file, Table S4). There was a little difference between the sites which we had identified as eligible to participate in the use-cases and those we had identified as not eligible (77% use-case sites *versus* 75% non-use-case sites).

Results micro-level data

The greater the number of coding systems in use, the harder it will be to achieve semantic interoperability; therefore, the

micro-level data collection was primarily concerned with collecting information about the coding systems the repositories used. We found that the WHO International Classification of Disease (ICD)²⁰ was the most common coding system used by 71% ($n = 39$) of respondents. ICD-10 ($n = 32$) was used by 82%; 13% ($n = 5$) used ICD-9; 23% ($n = 9$) used an ICD modification and 5% ($n = 2$) did not respond (Supplementary data file, Table S5).

The second most used coding system was the WHO International Classification of Primary Care (ICPC), this was used by 20% ($n = 11$) of respondents. Eighty-two percent ($n = 9$) of those using ICPC used ICPC-2 and 18% ($n = 2$) used ICPC-1, none reported using an ICPC modification (Supplementary data file, Table S6). The third most common coding system was the Systematised Nomenclature of Medicine (SNOMED),²¹ which was used by 13% ($n = 7$) of all the respondents; 44% ($n = 4$) of those using SNOMED used the Clinical Terms version; 33% ($n = 3$) used the Reference Terminology version and 22% ($n = 2$) did not respond (Supplementary data file, Table S7).

One of the least common coding systems used was the Read Coding system [version2 – 5-byte and the Clinical Terms Version 3 (CTv3)] and these were only used by the seven UK repositories. They represented 9% ($n = 5$) and 4% ($n = 2$) of all respondents, respectively; 87% ($n = 50$) did not respond (Supplementary data file, Table S8).

The survey highlighted that there was a great variety in the number of drugs dictionaries utilised by the repositories and this is one potential barrier to achieving semantic interoperability. Sixty percent ($n = 33$) of respondents said that they have a coding system for drugs (Primary care 83%, $n = 24$; Cancer 33%, $n = 5$; Genetic 36%, $n = 4$). Of these; 76% ($n = 25$) use the Anatomical Therapeutic Chemical classification system;²² 9% ($n = 3$) use Multilex; 12% ($n = 4$) responded 'other' and 3% ($n = 1$) responded 'no data' (Supplementary data file, Table S9). We were interested to know whether it was possible to extract information about the administration of drugs and we asked respondents if it was possible to extract data about daily dose and administration route from their database. Only around one-third of the Primary care and Cancer registries could extract data of this nature, while none of the genetic databases held this information (Supplementary data file, Table S10).

The survey was designed to assess what systems the registries had in place to achieve interoperability and to ensure data quality. Thirty-four percent ($n = 19$) of respondents had no system at all; only 5% ($n = 3$) used Health Level-7 (HL-7), an international interoperability organisation who's Reference Information Model underpins much interoperability in healthcare; 2% ($n = 1$) used the Clinical Data Interchange Standards Consortium (CDISC);²³ none used the Biomedical Research Integrated Domain Group (BRIDG)²⁴ and 52% ($n = 29$) used an 'In-house or other' system (Table 1). Nearly, all (93%, $n = 52$) of the respondents either had no system in place or used an in-house system or provided no data.

Table 1 Systems used to ensure data quality

	No system		HL-7		CDISC		BRIDG		In-house or other		No data	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Primary care	9	31	0	0	0	0	0	0	18	62	2	7
Cancer registry	4	25	2	13	1	6	0	0	7	44	2	12
Genetic 'Biobank'	6	55	1	9	0	0	0	0	4	36	0	0
Total	19	34	3	5	1	2	0	0	29	52	4	7

Data collection meso- and socio-cultural levels

Data extraction at this level was concerned with record level issues. The majority of respondents (82%; $n = 45$) have the ability to extract data in standardised formats such as Comma Separated Values, Excel and full text. All of the respondents have at least one appropriate format. The data collected have a wide application and this is reflected by the diverse nature of the information stored within these repositories which ranges from research to mortality records (Table 2).

The respondents reported that socio-cultural influences had a small but significant impact on the validity of their data. These factors included ethical, religious and legal factors (Table 3); these might delay or prevent participation in the TRANSFoRm studies.

Socio-cultural factors, which include legal and ethical constraints, as well as influences on diagnosis, and organisational components of the health system from which the data originates are often barriers to conducting research. In summary, 71% ($n = 39$) of respondents use ICD and 20% ($n = 11$) use ICPC; however, 86% ($n = 48$) do not use one of the three main systems for ensuring data quality; 29% opt instead for an in-house system. Very few sites are adopting national standards for interoperability in linking data. Whilst multiple drug dictionaries were used, 66% ($n = 10$) of cancer repositories did not use them. Extract formats for data were standardised and only 3% ($n = 6$) of respondents chose to use a non-standard format. Data were not forthcoming from eHR vendors ($n = 40$). Repositories had a broad range of applications for their data, the most important was research (49%, $n = 51$). The most common

socio-cultural influence that could potentially affect the validity of their data was ethical (10%, $n = 7$) and social (10%, $n = 7$) factors although 49% ($n = 36$) reported no social issues at all.

Difference in response depending on eligibility

Data sources that were non-use-case eligible tended to produce much fewer positive responses than those that were eligible. Overall, the repositories identified as potentially being use-case eligible made 2098 positive responses to questionnaire items compared with 268 from non-use-case eligible data sources (OR: 4.59; 95% CI: 3.93–5.35; $p < 0.008$); for genetic databases, the respective figures were 380:44 (OR: 6.13; 95% CI: 4.25–8.85; $p < 0.008$) and for cancer registries, they were 553:44 (OR: 5.87; 95% CI: 4.13–8.34; $p < 0.008$); the full results are in Table 4.

Data repositories capable of participation in the survey

Of the 56 valid responses, there were 15 pairs eligible to complete one or other of the use-cases. The 56 valid responses were made up of 29 databases of routine primary care data, 12 genetic databases and 15 cancer registries. From the valid responses, we were able to identify the location of databases with the potential to participate in the research studies. We identified five locations for linking primary care databases with genetic databases and 10 for linking primary care databases with cancer registries. The 15 eligible sites were spread across 11 countries (Supplementary data file, Table S11).

Table 2 The aims of the data source for the data collected

	Research		Mortality records		Financial monitoring		Quality performance monitoring		Other	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>N</i>	%	<i>N</i>	%
Primary care	27	43	5	8	8	13	14	23	8	13
Cancer registry	13	45	4	14	3	10	4	14	5	17
Genetic 'Biobank'	11	79	0	0	1	7	0	0	2	14
Total	51	49	9	9	12	11	18	17	15	14

Table 3 Socio-cultural influences on the validity of the data

	No social issues		Legal		Ethical		Social		Economic		Political		Religious		Other		No Data	
	n	%	n	%	n	%	N	%	n	%	n	%	n	%	n	%	n	%
Primary care	20	57	1	3	1	3	5	14	4	11	0	0	0	0	2	6	2	6
Cancer registry	9	43	3	14	3	14	0	0	2	10	0	0	1	5	1	5	2	9
Genetic 'Biobank'	7	40	1	6	3	18	2	12	0	0	1	6	2	12	1	6	0	0
Total	36	49	5	7	7	10	7	10	6	8	1	1	3	4	4	6	4	5

Table 4 Positive responses to the questionnaire sections – comparing non-use case eligible and use-case-eligible data sources

Level	Category	Use case eligible	Not use case eligible	OR	95% CI lower bound	95% CI upper bound	p
Primary care data sources							
MICRO	Data source	654	72	2.43	1.77	3.33	0.000
	Data interoperability	499	137	5.19	4.13	6.52	0.000
	Subtotal	1153	209	4.47	3.74	5.34	0.000
MESO	Record system	56	32	4.38	2.33	8.21	0.000
MACRO	Organisational	564	0				
	Socio cultural	174	6	14.50	5.91	35.59	0.000
	Subtotal	738	6	12.75	5.36	30.33	0.000
STUDY	Use-case specific	151	21	0.49	0.20	1.19	0.000
	Overall	2098	268	4.59	3.93	5.35	0.000
Genetic data/Biobanks							
MICRO	Data source	133	13	4.18	2.04	8.54	0.000
	Data interoperability	73	17	7.05	3.89	12.79	0.000
	Subtotal	206	30	7.32	4.70	11.39	0.000
MESO	Record system	16	5	7.47	1.39	40.25	0.000
MACRO	Organisational	106	0				
	Socio cultural	32	6	0.76	0.14	4.25	0.395
	Subtotal	138	6	0.63	0.12	3.20	0.001
STUDY	Use-case specific	20	24	2.22	0.52	9.51	0.007
	Overall	380	44	6.13	4.25	8.85	0.000
Cancer registries							
MICRO	Data source	210	5	12.86	4.89	33.82	0.000
	Data interoperability	128	23	7.97	4.77	13.33	0.000
	Subtotal	338	28	8.77	5.69	13.52	0.000
MESO	Record system	13	11	5.32	0.94	29.99	0.014
MACRO	Organisational	114	0				
	Socio cultural	45	5	3.00	0.88	10.18	0.000
	Subtotal	159	0				
STUDY	Use-case specific	43	0				
	Overall	553	44	5.87	4.13	8.34	0.000

Details of the eligible sites

The sites had a total of around 1.5 million potential patients eligible to participate in this research; over 30,000 in the genetics of diabetes use-case and over 1 million to participate in Barrett's disease, oesophageal cancer and the prescription of 30 medicines used to treat dyspepsia use-case. The country of origin, the website for these sites, the main coding system used and the expected delay in ethical approval are shown in Tables 6. We sometimes found contradictions between the data sources which indicated that they could supply linked data and several of the participants were, on closer questioning, only linking on a pilot basis; we have shaded out in grey the sites which are not currently active. The outcome of this process is that we have identified one fully functional location able to run the diabetes use-case (Table 5) and five pairs of locations able to run Barrett's disease use-case (Table 6). The one able to run the diabetes use-case is the Wellcome Type 2 Diabetes study group in Scotland. The five locations that can run the second use-case are as follows: Finland, Germany (Bremen), Norway, UK (General Practice Research Database), UK, Scotland (pilot).

DISCUSSION

Principal findings

The TIRRE survey has been completed by 56 data-repositories across Europe and six outside the EU. We have developed a usable instrument which can assess their potential to take part in linked data research. There were no equivalent International sites available to conduct this type of research. A challenge was to get databases to complete the questionnaire, when we did get a response, the completeness of information gathered was high and proved useful in identifying their potential to participate in linked research. Meso- and macro-level questions were important discriminators between use-case and non-use-case eligible data sources. There are currently no other survey instruments available to enable brokerage between databases potentially willing to participate in research. Micro levels informed about the data and its granularity.

Implications of the findings

The TIRRE survey is the first step towards assessing the potential of a database for linkage. It can identify data

Table 5 The eligible sites for conducting the diabetes TRANSFoRm use-cases

Pair	Name	Country	URL	Size (N/1000)	Date established	Primary coding system	Secondary coding system	Linkage Yes/ Pilot/ Plan/No	Ethics approval	Contact patients
									Months	Y/N
Primary Care & Genetic ‘Biobanks’										
1	Biobank of Medical University Graz	Austria	www.medunigraz.at/biobank	100–1000	2006	ICD-10	n/a	Pilot	≤3 months	Yes
	Austrian Primary Care Research Network	Austria	n/a	10–100	2005	ICPC-2	n/a	No	≤3 months	No
2	Généthon	France	www.genethon.fr	10–100	1991	n/a	n/a	Pilot	4–6 months	Yes
	OMG CONSTANCE GAZEL	France		0.1–1			n/a	No		
3	Da Vinci European Biobank (daVEB)	Italy	www.davincieuropeanbiobank.org	100–1000	2008	ICD-10	n/a	Yes	≤3 months	No
	Arianna database	Italy	n/a	10–100	2002	ICD-9	n/a	Pilot	n/a	Yes
4	SGI-RVB	Spain	www.csisp.gva.es/web/csisp	10–100	2010	SNOMED (CT)	Read Codes	Yes	≤3 months	Yes
	SIDIAP	Spain	n/a	>5000	2005	ICD-10	n/a	No	≤3 months	No
5	The Wellcome Trust Type 2 Diabetes Genetics Case-control Collection	UK-Scotland	www.diabetesgenetics.dundee.ac.uk	10–100	2004	n/a	n/a	Yes	≤3 months	Yes
	Health Informatics Centre	UK-Scotland	n/a	100–1000	1990	ICD-10	SNOMED	Yes	≤3 months	Yes

Table 6 The eligible sites for conducting Barrett's disease TRANSFoRm use-cases

Pair	Name	Country	URL	Size (N/1000)			Secondary coding system Y/N	Linkage Yes/ Pilot/ Plan/ No	Ethics approval	Contact patients
Primary care & cancer registry										
1	CroDiab	Croatia	www.idb.hr/web_english/crodiab.htm	100–1000	2000	ICD-10	n/a	No	No	No
	Croatian Institute for Health Insurance	Croatia	www.hzjz.hr	1–10	2001	ICD-10	ICPC-2	Yes	n/a	yes
2	Finnish Cancer Registry	Finland	www.cancerregistry.fi	1000–5000	1952	ICD-0-3	SNOMED	Yes	≤3 months	n/a
	Care Register for Health Care (HILMO) Institutions	Finland	www.stakes.fi	>5000	1961	ICD-10	ICPC-2	Yes	≤3 months	
3	Tumorotheque Régionale de Franche-Comté	France	www.chu-besancon.fr/tumoro	1–10	2005	ICD-10	n/a	Pilot	4–6 months	Yes
	OMG CONSTANCE GAZEL	France						No		
4	Bremen Cancer Registry	Germany	www.krebsregister.bremen.de	100–1000	1994	ICD-10	n/a	Pilot	≤3 months	Yes
	CONTENT	Germany	www.content-info.org	100–1000	2005	ICD-10	ICPC-2	Pilot	4–6 months	Yes
5	National Cancer Registry	Ireland	www.ncri.ie	n/a	1994	ICD-10	n/a	Pilot	No	Yes
	GP MED	Ireland	www.icgp.ie	100–1000	2008	ICPC-2	n/a	Pilot	≤3 months	No
6	Janus Serum Bank, Cancer Registry, Norway	Norway	www.kreftregisteret.no/en/	100–1000	1951	ICD-10	SNOMED (CT&RT)	Yes	≤3 months	Yes
	Norwegian Prescription Database	Norway	www.norpd.no	1000–5000	2004	ICD-10	ICPC-2	Yes	n/a	n/a
7	Spanish DILI Registry	Spain	www.spanishdili.uma.es	0.1–1	1994	n/a	n/a	Pilot	≤3 months	No
	SIDIAP	Spain	n/a	>5000	2005	ICD-10	n/a	No	≤3 months	No
8	Spanish Tumour Bank Network	Spain	www.cnio.es	100–1000	2000	ICD-10	n/a	Pilot	≤3 months	No
	SIDIAP	Spain	n/a	>5000	2005	ICD-10	n/a	No	≤3 months	No
9	National Cancer Intelligence Network (NCIN)	UK-England	www.ncin.org.uk/home.aspx	>5000	1971	ICD-10	n/a	Yes	4-6 months	Yes
	General Practice Research Database (GPRD)	UK-England	admin@gprd.com	>5000	1987	Read Codes	n/a	Yes	≤3 months	Yes
10	SCI-DC, Scotland	UK-Scotland		1000–5000	2005	ICD-10	Read Codes	Pilot	≤3 months	Yes
	Health Informatics Centre	UK-Scotland	n/a	100–1000	1990	ICD-10	SNOMED	Yes	≤3 months	Yes

sources suitability in terms of data availability and readiness to participate in a study. Whilst the initial focus of TIRRE was on linking data sources (which were important and consistent), the meso- and macro-factors generally had higher OR of predicting use-case eligibility.

Different coding systems have varying levels of granularity. For example, at the time of this study, neither ICD-10 nor ICPC differentiated between types of diabetes according to the latest WHO classification. ICD-10 differentiates insulin-dependent and non-insulin-dependent, rather than the Type 1 (insulin for survival) and Type 2 diabetes used in the latest classifications. Although we acknowledge, this is now updated in later releases.

Comparison with the literature

It is possible to draw comparisons between the complexity of this task and the existing successful projects that involve linking data. However, the successful data repositories in the UK have all been based on a single vendor of GP eHR system. Clinical Practice Research Datalink previously only extracted data from a single vendor called *In-Practice Systems*, though they are expanding this to all UK vendors;²⁵ Q-Research on the EMIS system²⁶ and other UK research networks (The Health Improvement Network²⁷ and ResearchOne²⁸) and other networks following the same pattern. The only exception to this in the UK is the Royal College of General Practitioners (RCGP) Research and Surveillance Centre (RSC);²⁹ this network extracts data from all the different brand of medical record systems. It has published a cohort profile about patients in the RCGP RSC database with diabetes, one of the TRANSFoRm use-case areas³⁰ Notwithstanding the RSCHP RSC success, the relatively simple task of linking data from this small number of brands of computer within the UK has proved challenging, both in terms of creating a summary care record³¹ and in developing a common data extraction system.³²

Limitations of the method

Any initial screening process will need to be followed up by a detailed assessment of whether the dataset needed for a given study can be elicited from the data repositories. There was no real incentive for data repositories to supply us with the data required, as there was not a reciprocal offer of benefit. As a consequence, our results inevitably underestimate the number of sites where this type of research can

be conducted. We propose that future projects should consider including incentives in their budget. An effective method to reduce the impact of this self-selection bias could be to approach databases with a partially completed survey (using information available in the public domain) in order to encourage participation. Furthermore, the collected data could be shared publicly as a metadata registry that would facilitate advertising data offered by organisations for prospective studies. We also recommend limiting surveys to 30–40 questions to improve the response rate.

Call for further research

We need to conduct test-retest studies to assess the reliability of the survey instrument. The reliability test could be carried out by repeating the data collection after a period of time. While this would help to validate the instrument, it will also potentially remove any bias introduced by the specific person responding to the survey. We should conduct simulated and real studies with data extractions to test its validity. However, conducting real studies may be affected by the availability of funding. Alternatively, we can promote reuse of the instrument in other projects with the research area.

CONCLUSIONS

A large complex set of data is needed to know if it will be possible to link primary care and either disease registry of the genetic database. This complex set of data can either be classified by level of granularity or as a business or data requirement.

The TIRRE instrument is a useful tool that can be used to assess general suitability and readiness to participate in linked research studies. With increased use, it is likely that TIRRE will evolve further, but its use needs to be embedded in a concrete 'offer' and business case rather than a one-off research study.

Acknowledgements

Paul van Royen for his comments on the manuscript; IMIA and EFMI for supporting their primary health care informatics working groups. Antonis Ntasioudis for his contribution to this research. TRANSFoRm is supported by the European Commission – DG INFSO (FP7 2477).

REFERENCES

1. Trifirò G, Coloma PM, Rijnbeek PR, Romio S, Mosseveld B, Weibel D, *et al.* Combining multiple healthcare databases for postmarketing drug and vaccine safety surveillance: why and how? *Journal of Internal Medicine* 2014;275(6):551–61.
2. Bohensky MA, Jolley D, Sundararajan V, Evans S, Pilcher DV, Scott I, *et al.* Data linkage: a powerful research tool with potential problems. *BMC Health Services Research* 2010;10:346.
3. Boyd JH, Randall SM, Ferrante AM, Bauer JK, Brown AP, Semmens JB. Technical challenges of providing record linkage services for research. *BMC Medical Informatics and Decision Making* 2014;14:23.
4. van Panhuis WG, Paul P, Emerson C, Grefenstette J, Wilder R, Herbst AJ, *et al.* A systematic review of barriers to data sharing in public health. *BMC Public Health* 2014;14:1144.
5. de Lusignan S, Chan T. The development of primary care information technology in the United Kingdom. *The Journal of Ambulatory Care Management* 2008;31(3):201–10.

6. de Lusignan S, van Weel C. The use of routinely collected computer data for research in primary care: opportunities and challenges. *Family Practice* 2006;23(2):253–63.
7. Kellerman AL, Jones SS. What will it take to achieve the as-yet-unfulfilled promise of health information technology? *Health Affairs* 2013;32(1):63–8.
8. Delaney BC, Curcin V, Andreasson A, Arvanitis TN, Bastiaens H, Corrigan D, et al. Translational medicine and patient safety in Europe: TRANSFoRm—architecture for the learning health system in Europe. *BioMed Research International* 2015;2015:961526. doi: 10.1155/2015/961526
9. Mastellos N, Bliźniuk G, Czopnik D, McGilchrist M, Misiaszek A, Bródka P, et al. Feasibility and acceptability of TRANSFoRm to improve clinical trial recruitment in primary care. *Family Practice* 2016;33(2):186–91.
10. Ethier JF, Curcin V, McGilchrist MM, Keung SNLC, Zhao L, Andreasson A, et al. eSource for clinical trials: implementation and evaluation of a standards-based approach in a real world trial. *JMIR* 2017;106:17–24.
11. European Union; Europe's Information Society. Communication from the Commission to the Council, the European Parliament, the European Economic and Social Committee and the Committee of the Regions-e-Health – making healthcare better for European citizens: an action plan for a European e-Health Area {SEC(2004)539}; 2004. Available from: <http://eurlex.europa.eu/LexUriServ/LexUriServ.do?uri=CELEX:52004DC0356:EN:NOT>
12. de Lusignan S, Krause P, Michalakidis G, Vicente MT, Thompson S, McGilchrist M, et al. Business process modelling is an essential part of a requirements analysis. Contribution of EFMI Primary Care Working Group. *Yearbook of Medical Informatics* 2012;7(1):34–43.
13. Leppenwell E, de Lusignan S, Vicente MT, Michalakidis G, Krause P, Thompson S, et al. Developing a survey instrument to assess the readiness of primary care data, genetic and disease registries to conduct linked research: TRANSFoRm International Research Readiness (TIRRE) survey instrument. *Informatics in Primary Care* 2012;20(3):207–16.
14. de Lusignan S, Cashman J, Poh N, Michalakidis G, Mason A, Desombre T, et al. Conducting requirements analyses for research using routinely collected health data: a model driven approach. *Studies in Health Technology and Informatics* 2012;180:1105–7.
15. de Lusignan S, Pearce C, Shaw NT, Liaw ST, Michalakidis G, Vicente MT, Bainbridge M; International (IMIA); European (EFMI) Medical Informatics Association; Federation Primary Care Informatics Working Groups (PCI-WG). What are the barriers to conducting international research using routinely collected primary care data? *Studies in Health Technology and Informatics* 2011;165:135–40.
16. de Lusignan S, Liaw ST, Krause P, Curcin V, Vicente MT, Michalakidis G, et al. Key concepts to assess the readiness of data for international research: data quality, lineage and provenance, extraction and processing errors, traceability, and curation. Contribution of the IMIA Primary Health Care Informatics Working Group. *Yearbook of Medical Informatics* 2011;6:112–20.
17. Dolin RH, Alschuler L. Approaching semantic interoperability in Health Level Seven. *Journal of the American Medical Informatics Association* 2011;18(1):99–103.
18. Michalakidis G, Kumarapeli P, Ring A, van Vlymen J, Krause P, de Lusignan S. A system for solution-orientated reporting of errors associated with the extraction of routinely collected clinical data for research. *Studies in Health Technology and Informatics* 2010;160:724–8.
19. Santos M, Bax M, Kalra D. Building a logical EHR architecture based on ISO 13606 standard and semantic web technologies. *Studies in Health Technology and Informatics* 2010;160(Pt 1):161–5.
20. International Classification of Diseases. Available from: <http://www.who.int/classifications/icd/en/>. Accessed 15 December 2018.
21. NHS Digital. Terminology and classifications. Available from: <https://digital.nhs.uk/article/290/Terminology-and-Classifications>. Accessed 15 December 2018.
22. Anatomical Therapeutic Chemical classification system (ATC). Available from: http://www.whocc.no/atc/structure_and_principles/. Accessed 15 December 2018.
23. Clinical Data Interchange Standards Consortium (CDISC). Available from: <http://www.cdisc.org>. Accessed 15 December 2018.
24. Biomedical Research Integrated Domain Group (BRIDG). Available from: <http://www.cdisc.org/bridg>. Accessed 15 December 2018.
25. Kousoulis AA, Rafi I, de Lusignan S. The CPRD and the RCGP: building on research success by enhancing benefits for patients and practices. *British Journal of General Practice* 2015;65(631):54–5.
26. Hippiusley-Cox J, Stables D, Pringle M. QRESEARCH: a new general practice database for research. *Informatics in Primary Care* 2004;12(1):49–50.
27. Blak BT, Thompson M, Dattani H, Bourke A. Generalisability of The Health Improvement Network (THIN) database: demographics, chronic disease prevalence and mortality rates. *Informatics in Primary Care* 2011;19(4):251–5.
28. ResearchOne Health and Care Database. Available from: <http://www.researchone.org>. Accessed 4 August 2018.
29. Correa A, Hinton W, McGovern A, van Vlymen J, Yonova I, Jones S, et al. Royal College of General Practitioners Research and Surveillance Centre (RCGP RSC) sentinel network: a cohort profile. *BMJ Open* 2016;6(4):e011092.
30. McGovern A, Hinton W, Correa A, Munro N, Whyte M, de Lusignan S. Real-world evidence studies into treatment adherence, thresholds for intervention and disparities in treatment in people with type 2 diabetes in the UK. *BMJ Open* 2016;6(11):e012801.
31. Greenhalgh T, Stramer K, Bratan T, Byrne E, Russell J, Potts HW. Adoption and non-adoption of a shared electronic summary record in England: a mixed-method case study. *BMJ* 2010;340:c3111.
32. NHS Digital. GP Extraction Service. Available from: <http://content.digital.nhs.uk/gpes>. Accessed 15 December 2018.

Appendix 1 Details of the TRANSFoRm work tasks

Description of work tasks (WT) 6.1 and 6.2 of the TRANSFoRm project:	
WT 6.1:	Requirements analysis of EHRs
1.	Using the results of the EHR capacity study within the EGPRN and ESPCG networks (WT1.1), we will conduct an in-depth study of the most common EHRs systems used in Europe to examine the availability of API details. The scope of the study will include patient-held records, which may hold substantially less coded and structured data. These will include Microsoft and Google patient record systems – and countries where health cards are used.
2.	We will conduct a parallel in-depth study of data repositories that can be used for clinical trials. We will look to identify local, EHR brand specific and health-system access points to primary care data. The types of data-access points we might be able to run queries on might include: 1) Billing or performance indicator extracts of routine data; 2) Sentinel networks or research network database and 3) National data extract systems with closed API. These may provide pragmatic quick win access to primary care data while a longer-term access is being developed.
WT 6.2:	Requirements analysis of genotype data repositories.
	We will conduct a parallel in-depth study of (genotype) clinical data repositories across Europe and their potential use for clinical research. The scope of the study will include structured genotype data and potential integration points with patient healthcare for biomedical and translational clinical research. Genotype data is normally held by Biobanks or other research organisations either as sample identification information or specific codes for Single Nucleotide Polymorphisms (SNPs).

Table S1 Categories of data collection and min-to-max number of questions; skip logic reduces the number of questions that each type of respondent might answer

Level	Category	Primary care data	Genetic database	Cancer registry	Others (social care data, cohorts)	Number of questions
Micro	Data source	30–46	36–54	30–46	33–51	54
	Data interoperability	30–43	31–43	31–43	31–43	43
	Subtotal	60–89	67–97	61–89	64–94	97
Meso	Record system	5–30				30
Macro	Organisational	15–15				15
	Socio-cultural	6–10				10
	Subtotal	21–25				25
Study	Use-case specific	5–8				8
Total		Min: 91	Min: 98	Min: 92	Min: 95	160
		Max: 152	Max: 160	Max: 152	Max: 157	

Table S2 Number of contacts and valid responses

On-line completion of the TIRRE questionnaire				
Online	Primary care data	Genetic data	Cancer registry data	Total (n)
Partial	21	11	13	45
Valid	12	9	10	31
Total	33	20	23	86

Table S3 Number of contacts and valid responses

Number of contacts and valid responses				
	Primary care (n)	Genetic (n)	Cancer registry (n)	Total (n)
First contacted by email	110	117	89	316
First contacted by phone	83	19	16	118
Other phone contacts	53	25	17	95
Declined	1	1	0	2
Partial responses	39	14	18	71
Valid responses	29	12	15	56

Table S4 Completion of the questionnaire

Analysis by respondent	Number of questions								
	Potential maximum responses	Primary care		Cancer registry		Genetic		Overall	
		Mean	Response	Mean	Response	Mean	Response	Mean	Response
Categorical questions	<i>n</i>	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Data source	35	22.41	64	24.00	69	25.00	71	23.39	67
Data interoperability	31	29.34	95	28.80	93	26.25	85	28.54	92
Micro-level	66	51.76	78	52.80	80	51.25	78	51.93	79
Meso-level/record system	12	3.10	26	2.33	19	2.67	22	2.80	23
Organisational	15	15.00	100	15.00	100	15.0	92	14.73	98
Socio-cultural	6	5.17	86	5.47	91	4.50	75	5.11	85
Macro-level	21	20.17	96	20.47	97	18.25	87	19.84	94
Overall	104	79.90	76.82	80.60	77.50	83.27	80.07	80.76	77.66

Table S5 Coding systems information (ICD)

	Is ICD used?						Version						Modified			
	Yes		No		Don't Know		ICD-9		ICD-10		No Data		Yes (number)		No/no data	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Primary care	20	69	9	31	0	0	3	15	17	85	0	0	5	25	15	75
Cancer registry	12	80	2	13	0	0	1	8	9	75	2	5	3	25	9	75
Genetic 'Biobank'	7	64	3	27	1	4	1	14	6	86	0	0	1	14	6	86
Total	39	71	14	25	1	4	5	13	32	82	2	5	9	23	30	77

Table S6 Coding systems information (ICPC)

	Is ICPC used?						Version				Modified			
	Yes		No		Don't Know		ICPC-1		ICPC-2		Yes		No	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Primary care	11	38	17	17	1	3	2	18	9	82	0	0	11	100
Cancer Registry	0	0	12	12	3	20	0	n/a	0	n/a	0	0	0	0
Genetic 'Biobank'	0	0	10	10	1	9	0	n/a	0	n/a	0	0	0	0
Total	11	20	39	71	5	9	2	18	9	82	0	0	11	100

Table S7 Coding systems information (SNOMED)

	Is SNOMED used?						Version					
	Yes		No		Don't Know		CT		RT		No Data	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Primary care	3	10	24	83	2	7	1	33	1	33	1	34
Cancer registry	2	13	11	73	2	14	1	33	1	33	1	34
Genetic 'Biobank'	2	18	7	64	2	18	2	67	1	33	0	0
Total	7	13	42	76	6	11	4	44	3	33	2	23

Table S8 Reading coding systems usage (CTv3 and Read codes version used)

	Are the Clinical Terms version 3 used (CTv3)?		Are Read Codes used?		None		No Data	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Primary care	2	6	3	10	25	81	1	3
Cancer registry	0	0	1	7	11	73	3	20
Genetic 'Biobank'	0	0	1	9	10	91	0	0
Total	2	4	5	9	46	81	4	6

Table S9 Coding systems for drugs

	Is there a coding system for drugs?						What coding system									
	Yes		No		Don't Know		ATC Codes		DAAD		Multilex		Other		No Data	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Primary care	24	83	65	17	0	0	19	79	0	0	2	8	3	13	0	0
Cancer registry	5	33	9	60	1	7	3	60	0	0	1	20	1	20	0	0
Genetic 'Biobank'	4	36	6	55	1	9	3	75	0	0	0	0	0	0	1	25
Total	33	60	20	36	2	4	25	76	0	0	3	9	4	12	1	3

Table S10 Extraction of drug information from data provided

	Daily dose		Frequency		Administration route		No Data	
	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%	<i>n</i>	%
Primary care	20	35	18	32	19	33	0	0
Cancer registry	3	3	3	30	2	20	2	20
Genetic 'Biobank'	0	0	0	0	0	0	4	100
Total	23	32	21	30	21	30	6	8

Table S11 Location of respondents and eligible sites

Countries	EU	Routine primary care data		Genetic 'Biobank'		Cancer registry		Use-case location	
		<i>N</i>	<i>n</i>	<i>N</i>	<i>n</i>	<i>N</i>	<i>n</i>	Primary Care + Genetic/Biobank	Primary Care + Cancer Registry
Austria	X	5	1	3	1	4	0	X	
Belgium	X	10	2	3	0	3	0		
Bulgaria	X	0	0	0	0	0	0		
Croatia		1	1	0	0	1	1		X
Cyprus	X	0	0	0	0	0	1		
Czech Republic	X	0	0	0	0	0	2		
Denmark	X	4	1	1	0	1	0		
Estonia	X	3	0	0	1	1	0		
Finland	X	1	1	1	0	1	1		X
France	X	11	3	24	1	13	1	X	X
Germany	X	9	4	18	0	4	1		X
Greece	X	4	0	1	0	1	1		

Table S11 (Continued)

Countries	EU	Routine primary care data		Genetic 'Biobank'		Cancer registry		Use-case location	
		N	n	N	n	N	n	Primary Care + Genetic/Biobank	Primary Care + Genetic/Biobank
Hungary	X	2	0	0	0	0	0		
Iceland		2	0	5	1	2	0		
Ireland	X	4	1	1	0	1	1		X
Italy	X	15	3	27	1	18	0	X	
Latvia	X	0	0	1	1	1	0		
Lithuania	X	1	0	1	1	1	0		
Luxembourg	X	1	0	0	0	0	0		
Malta	X	1	1	3	0	2	0		
Norway		5	1	2	0	2	1		X
Poland	X	0	0	1	0	4	1		
Portugal	X	6	0	2	1	3	1		
Romania	X	7	0	2	0	1	0		
Russia		2	0	0	0	1	0		
Slovakia	X	1	1	0	0	1	0		
Slovenia	X	0	0	0	0	1	0		
Spain	X	17	1	11	1	14	2	X	XX
Sweden	X	7	1	1	0	1	0		
Switzerland		1	1	2	1	9	0		
The Netherlands	X	24	2	6	0	4	0		
Turkey		2	0	0	0	0	0		
UK-England	X	11	3	3	0	11	1		X
UK-Northern Ireland	X	0	0	0	0	1	0		
UK-Scotland	X	6	1	2	1	1	1	X	X
UK-Wales	X	0	0	0	0	0	0		
European*	X	2	0	9	1	1	1		
Total		165	29	130	12	111	15	5	10