

BMJ Open

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

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

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

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

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

BMJ Open

Existing Barriers and Recommendations of Real-World Data Standardization for Clinical Research in China: A Qualitative Study

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-059029
Article Type:	Original research
Date Submitted by the Author:	09-Nov-2021
Complete List of Authors:	Lai, Junkai; Peking University First Hospital, Peking University Clinical Research Institute, Liao, Xiwen; Peking University First Hospital, Department of Biostatistics Jin, Feifei; Peking University People's Hospital, Department of Trauma Medicine Wang, Bin; Peking University First Hospital, Department of Biostatistics Yao, Chen; Peking University First Hospital, Department of Biostatistics; Peking University Clinical Research Institute Li, Chen; Fourth Military Medical University, Department of Health Statistics; School of Preventive Medicine
Keywords:	QUALITATIVE RESEARCH, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

1
2
3 **Existing Barriers and Recommendations of Real-World Data Standardization for**
4 **Clinical Research in China: A Qualitative Study**
5
6

7 Junkai Lai¹, Xiwen Liao¹, Feifei Jin², Bin Wang¹, Chen Yao^{1,3,4}, Chen Li⁵
8
9

10 1 Department of Biostatistics, Peking University First Hospital, Beijing, Beijing, China.

11 2 Peking University People's Hospital, Beijing, Beijing, China.

12 3 Peking University Clinical Research Institute, Beijing, Beijing, China.

13 4 Hainan Institute of Real World Data, Qionghai, Hainan, China.

14 5 Department of Health Statistics, School of Preventive Medicine, Fourth Military
15 University, Xi'an, Shaanxi, China.
16
17

18
19
20 **Correspondence to**

21 Chen Yao

22 Peking University First Hospital, Xicheng District, Beijing 100034, China

23 Tel +86 18610640562

24 Email yaochen@hsc.pku.edu.cn
25
26
27

28 **Word Count**

29 3996
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Existing Barriers and Recommendations of Real-World Data Standardization for Clinical Research in China: A Qualitative Study

Abstract

Objective

To investigate the existing barriers and recommendations of real-world data standardization for clinical research through a qualitative study on different stakeholders.

Design

This qualitative study involved five types of stakeholders based on five interview outlines. The data analysis was performed using the constructivist grounded theory analysis process.

Setting

8 Hospitals, 4 Hospital System Vendors, 3 Big Data Companies, 6 Pharmaceutical Companies, and 4 Regulatory Institutions were included.

Participants

In total, 62 participants from 25 institutions were interviewed through purposive sampling.

Results

Real world data is difficult to standardize for clinical research in China. The main causes were difficulty integrating standards in routine clinical process, lack of generalizable research datasets, and difficulty auditing the standardization process. The main suggestions were better feedback loop for standards, better separation of standards and documentation, promoting generalizable clinical research data models, and improving traceability to source data.

Conclusions

Determining the barriers and recommendations could contribute to the development of clinical research in China. The key findings in the study make a clear point that data standardization relies on the consensus of many working groups. Data standards cannot be integrated at the source data unless it is customized for practical usage and can only be secondarily applied to source data adequately when the meaning of the source data can be clearly understood.

Keywords: real world data; data standards; data standardization; clinical research; qualitative research; China;

Strengths and Limitations of this study:

- Strength: Wide variety of relevant stakeholders on the subject
- Strength: Qualitative understanding of a major industry bottleneck

- Strength: Important recommendations that can guide the direction of the future of the subject
- Limitations: Due to COVID-19, a portion of the interviews were not done in person and may limit the ability to read into the participants response for further exploration of the subject
- Limitations: Recruitment of participants were limited by those that were already exploring the subject and cannot be generalized to participants that may not be familiar with the subject.

Introduction

Real World Data are data relating to patient health status or the delivery of health care routinely collected from a variety of sources.¹⁻⁴ Data standards related to real-world data sources (such as EMR data in China) have gradually evolved from the basic guidelines of clinical documentation to standard terminology usage and clinical data models.⁵⁻⁷ Data standard usage is meant to improve the data quality of real-world data so that it can be further used for secondary purposes such as clinical decision support and clinical research.

Since 2008, the Informatization Leading Group of the Ministry of Health and other departments have organized relevant experts to write the basic architecture and data standard specification of the Electronic Medical Record (EMR).⁶ By considering the basic guides for case writing that integrates Chinese and Western medicine and a survey on the common medical functions, the institution published a guidance for basic standards associated with EMRs.⁷ Through several national strategies to promote EMR meaningful use and data sharing, China has gradually improved the standardization of real-world data.

In 2018, a performance rubric for the meaningful usage of EMRs were created to rate hospitals;⁸ The policy pushed for all tertiary hospitals to reach grade 4 or above (hospital wide information sharing and primary medical decision support) and secondary hospitals to reach grade 3 or above (inter-departmental data exchange) by 2020 which has been overall achieved. Later in 2020 when EMRs have matured at hospitals, the Statistical Information Center of the National Health Commission further pushed for the data sharing capability of hospitals through another performance rubric.⁹ The main guidance on data standards currently is based on a localized version of HL7 CDA (Health Level 7 Clinical Document Architecture) and international terminology standards such as ICD-10 (International Classification of Disease version 10).⁶⁻⁷

The future of clinical research will benefit greatly if EMR data is standardized and able to be shared. However, even if there are existing guidance on industry health data standards, it is unclear what the barriers or best methods are for real-world data standardization for clinical research. Drawing on our previous qualitative study of the gap between real world data and clinical research, data standards are one of the components that causes this gap.¹⁰ This study seeks to address this question through a qualitative study by interviewing different stakeholders seeking to conduct clinical research using real-world data to benefit their decision-making process.

Methods

Design

Qualitative research allows us to understand a participants experience through qualitative methods of capturing data often through interviews. Constructivist grounded theory (CGT) provides a way for theory construction from qualitative data and is of the view that researchers do not discover theory but rather construct it through interaction and interpretation of the participant.¹¹⁻¹² Exploring why real-world data is difficult to standardize for usage in clinical research is the concern of this study. Therefore, a qualitative research strategy guided by CGT was employed.

The research team conducted in-depth interviews with different institutions representing 3 categories of key stakeholders: stakeholders affecting the source data, stakeholders affecting the standardization of source data for research, and stakeholders affecting the validity of data for clinical research. The interviews were conducted between August and October 2021. The study is reported following the guidelines of the Consolidated Criteria for Reporting Qualitative Research (COREQ) guidelines.¹³

Participants Selection

We intended to interview three categories of stakeholders: stakeholders affecting the source data, stakeholders affecting the standardization of source data for research, and stakeholders affecting the validity of data for clinical research. Hospital and Hospital System Vendors represented the first category, Big Data Companies represented the second category, and Pharmaceutical and Regulatory Institutions represented the third category. Purposive sampling was used for the selection of participants which represented institutions that had a team working with real world data for clinical research and recommended experts.¹⁴⁻¹⁵ The participants recruited for the study included 25 institutions with a total of 62 participants with no refusal in participation or dropped out. YC and JL contacted the interviewee through phone and briefed the subject matter and the objective of investigation before the participants agreed to be arranged for an interview. Interviewees represents their own opinions based on their experience working at the institution and do not represent the institution. The number of participants by the type of stakeholder is shown in shown in **Table 1**. Detailed list of institutions by type of stakeholder is included in the Appendix.

The inclusion criteria for the interviewees were as follows

Inclusion criteria

1. Participant had extensive experience as a staff member at stakeholder's institution
2. Participant has experience evaluating real world data for clinical research for the institution

Exclusion criteria

1. Participant could not sign informed consent form
2. Participant could not provide at least 45 minutes for an interview

Setting

1
2
3 The research team with training and experience in qualitative methods conducted
4 interviews using a telephone or in person. A quiet meeting room was chosen for each
5 interview to allow for better recording of the study data and exclusion of non-
6 participants.
7
8
9

10 **Data Collection**

11 Semi-structured interviews were recorded either over the phone or in person through
12 an iPhone app with the ability to transcribe audio into text files.¹⁶⁻¹⁸ Field notes were
13 taken to summarize important findings during the interview process that would help
14 guide later coding. Depending on the stakeholder's time and willingness, a focus
15 group interview was arranged instead of one-on-one interviews to promote discussion
16 and communication.¹⁹ The allowed time for interviews must not exceed 60 minutes in
17 respect to the participants daily schedule. After reading a confidentiality and privacy
18 statement and informing participants that interview will be recorded, the researcher will
19 conduct interviews in accordance with steps already described in this protocol. In each
20 interview, basic information of interview time, place and interviewee will be collected.
21 Five sets of interview guides were designed for the five types of stakeholders of the
22 interviewees and pilot tested before hand with similar participants not included in the
23 study to make the flow of questioning better. Full interview guides are included in the
24 Appendix. The interview questions guide the interviewer in exploring the subject with
25 the participant. Further discussion on the questions or repeat interviews were allowed
26 to explored deeper into the topic or for better clarification but in the study none occurred.
27 Transcripts were not returned to the participants for correction. The collection of data
28 was not limited by data saturation and finished once all participants were interviewed.

29 The interviewers were four doctoral students. JL (Male) and XL (Female) were
30 mainly responsible for the interviews, and BW (Male) and FJ (Female) played
31 supportive roles and were mainly responsible for the recordings. The interviewers were
32 trained in a qualitative research course and had experience conducting interviews.
33
34
35
36
37
38
39
40

41 **Sample Interview Guideline for the Hospitals**

- 42 ➤ This interview will be recorded and used in a qualitative study, your identity will
43 be concealed to protect your privacy, do we have your full consent in this interview
44 and have your signed information consent form?
- 45 ➤ Describe your role and experience facilitating clinical research at the Hospital?
- 46 ➤ What are the motivating goals of clinical research?
- 47 ➤ How do you determine the research dataset that are needed for your research?
- 48 ➤ Do you think that electronic medical records or routine care data at the hospital are
49 enough to accomplish your research?
- 50 ➤ How do you aggregate and store all data from different hospital systems?
- 51 ➤ How do you implement standards on your data?
- 52 ➤ What areas in your clinical research process do you have to rely on external vendors
53 to help you?
- 54 ➤ How does data sharing happen for medical records inside and outside of the
55 hospital?
56
57
58
59
60

- 1
2
3 ➤ Beyond clinical research, have you standardized your data for other purposes?
4
5

6 **Analysis**

7 All interviews were transcribed to text using the automated transcription software and
8 double checked for each recording by the two interviewers (JL and XL). Coding and
9 memoing were done by the four researchers (JL, XL, FJ, BW) whom drew on the
10 techniques of constructivist grounded theory while analyzing the data.¹⁶⁻¹⁸ QSR NVivo
11 V.12 software was used for coding. The team developed a structured coding tree based
12 on the interviews that started with inductive open coding. Once the core categories
13 emerged, deductive selective coding was performed. Open coding was performed
14 independently by the two researchers, and the derived core categories were compared
15 in multiple rounds of discussions until all four research members (JL, XL, FJ, BW)
16 agreed. Participants did not provide feedback on the findings.
17
18
19
20

21 **Patient and public involvement**

22 There was no patient or public involvement in this research.
23
24

25 **Results**

26 **Barrier and Suggestions in Data Standardization of Real-World Data for Clinical** 27 **Research**

28 The CGT framework generated from the three stages of coding and the 62 participants'
29 responses are summarized in the flow chart (figure 1). We found three causes that create
30 the barrier in data standardization of real-world data, including difficulty integrating
31 standards with routine clinical process, lack of generalizable research datasets, and
32 difficulty auditing the standardization process. The main suggestions relating to the
33 causes were found to be better feedback loop for standards, separation of data standards
34 and documentation, promoting clinical research data models, and improving
35 traceability to source data.
36
37
38
39
40

41 **Causes**

42 *Hard to Integrate Standards with the Routine Clinical Process*

43 Standard terminology libraries do not map well to the expressions used by physicians
44 to describe the patient's conditions even when translated since they do not incorporate
45 easily searchable colloquial terms. Standardized lists given as options will often lead to
46 physicians using the "other" option to fill in their answers. Data standards are rarely
47 utilized for the communication of data within hospitals and only standardized later on
48 for research or regulatory data submission purposes; In addition, there are lots of
49 isolated systems that may not be able to communicate with each other. Even if hospital
50 system vendors want to use default standards incorporated into the system, hospitals
51 will often reject them in favor of what they are familiar with. For research, specialty
52 departments with hospitals rely on external vendors to create databases for them with
53 standardized data; Physicians find it dissatisfying to give a search criterion and retrieve
54 datasets that find patients that do not match their initial inclusion criteria due to the
55
56
57
58
59
60

1
2
3 granularity of the standards.
4
5

6 “We give our clients default standards to use but they may feel that the standards do
7 not match their needs and will ask us to perform more customizations” – Hospital
8 System Participant 1
9

10
11 “When implementing standard answers for the diagnosis field, doctors often just fill in
12 their own answers in the “other” option” – Hospital Participant 8
13
14

15 *Lack of Generalizable Research Dataset*

16 Pharmaceutical companies hoping to use real world data as part of their product
17 development process find it to be a struggle. The usage of existing real-world data is
18 often limited to lab and demographic data that may provide limited use. The developed
19 hospital disease specialty databases by different big data vendors may differ greatly in
20 terms of data definition and cannot be externally validated or aggregated to conduct
21 studies. Regulatory institutions also express that the data shared by hospitals from data
22 sharing policies are limited in usage and cannot be directly applied for decision making.
23 Real world data currently lacks generalized data models that can be used for different
24 therapeutic areas as well as regulatory evaluations.
25
26
27
28

29 “For feasibility studies, we may look into disease specialty databases. The data
30 elements in these databases are usually very different from each other and we may have
31 to focus on data elements that are more widely available to conduct our studies” -
32 Pharmaceutical Participant 7
33
34
35

36 “Beside our department, other departments are also using real world data. Although
37 data quality may not be great, there may still be important signals that can support the
38 evaluation process. Developing a platform that can be used by multiple departments
39 may be in our interest.” – Regulatory Participant 3
40
41

42 *Hard to Audit Data Standardization Process*

43 Using existing databases, pharmaceutical companies may find it difficult to evaluate
44 the quality of data based on plausibility measures alone. Without a way to guarantee
45 traceability of data back to the source data, real world data is limited to conducting
46 feasibility studies and hard to incorporate into clinical studies such using real world
47 data as an external control group. Pharmaceutical companies also notice the usage of
48 advanced algorithms for the transformation of real-world data to research datasets and
49 find it concerning. The main concern is that new methodologies that risk data integrity
50 may be rejected by regulatory institutions during evaluation. Big data companies find
51 that advanced algorithms may produce unexpected and hard to explain variations in
52 new data sources and are striving to reduce the complexity in data transformation.
53
54
55
56
57

58 “From our experience, our team has had to reduce the amount steps for the
59 transformation of data from one model to another model and even resorted to merging
60

1
2
3 groups of IT staff working on different problems to clear up the standardization process”
4 – Big Data Participant 5
5
6

7 “My concerns for the usage of artificial intelligence algorithms for the extraction and
8 standardization of data are whether regulatory institutions will accept them.” –
9 Pharmaceutical Participant 4
10
11

12 **Suggestions**

13

14 *Better Feedback Loop for Standards*

15 Employees of hospital systems vendors will attend workshops and certification
16 programs by HL7 China to help implement more standardized data communication
17 between hospital systems. They find localization and extension of international
18 standards to be the primary focus in their efforts and require a joint effort between
19 participants from different institutions guided by standards organizations. Big data
20 companies and hospital system vendors also have rich experience customizing hospital
21 systems and standardizing local data and feel that they can enrich the development of a
22 more suitable local standard.
23
24
25
26
27

28 “When working to develop different research databases, our team has incorporated
29 medical experts that help us aggregate terminology libraries and understand the best
30 way to search for data.” – Big Data Participant 15
31
32

33 “Standards will get adopted if they make our systems communicate better with each
34 other and services well-liked by our clients.” – Hospital System Participant 4
35
36

37 *Better Separation of Data Standards and Documentation*

38 Good balance in separation of data standardization and routine documentation of data
39 is an important strategy to reduce the burden of physicians during documentation.
40 Hospitals will negotiate with vendors on mapping locally used terminology and
41 standard terminologies in the backend of the system to alleviate differences while
42 ensuring data can be standardized. In addition, hospital specialty departments may
43 choose to implement recommended text or other methods that simplify documentation
44 while enhancing consistency in data capture. Finally, hospitals will usually consult
45 external vendors that can leverage technology such as natural language processing that
46 deals better with standardizing electronic text data.
47
48
49
50

51 “Sometimes it is necessary for there to be some standardization during data collection
52 because doctors who are busy will often elaborate very little about the patient. For some
53 specialty departments, we may implement text recommendations to help standardize
54 the documentation” – Hospital Participant 9
55
56

57 “Doctors are unfamiliar with the different standards. We will usually work with
58 companies that can use better technology to help us standardize the data.” – Hospital
59
60

Participant 13

Promoting Clinical Research Data Models

Data is often needed for many projects or services in a given company and may benefit from having a data model that will reduce the repetition of work done to gather data. Big data and hospital system vendors have been developing their own data models inspired by approaches from international data models such as HL7, OHDSI (Observational Health Data Sciences and Informatics), CDISC (Clinical Data Interchange Standards Consortium), and other organizations to meet this goal. As more research is done using real world data, big companies have been able to gather core datasets needed for different types of clinical studies that they can use to increase the generalizability of their data model.

“Learning from Huawei’s and Alibaba’s approach to organize their services, we are starting to apply the HL7 RIM (Health Level 7 Reference Information Model) model to build a middle layer in which our different hospital system can create their services. Eventually we would like to use it to support clinical decision support systems” – Hospital System Participant 1

“When we participate in more clinical studies, we find better overlap between our schema and their research case report forms. Sometimes a therapeutic area may be very specific and we may need to extend our data model.” – Big Data Participant 5

Improving Traceability to Source Data

Regulatory institutions recommend that real world data standardization for clinical research should adhere to GCP (Good Clinical Practice) principles of data integrity in which data traceability is a key focus. Furthermore, they suggest that better integration between the collection of real-world data and clinical trial data management processes will lead to better regulatory acceptance. Prompted by the concerns for data traceability, pharmaceutical companies are exploring methods of data capture that can meet regulatory expectation while reducing data collection efforts.

“GCP principles should be upheld similarly when using real world data for clinical research. Applying aspects of the clinical trial workflow may be needed to raise the confidence in real world data collection.” – Regulatory Institution Participant 2

“We have been searching for eSource capability that can help us collect reliable data that can be easily audited and used as evidence for regulatory approval” - Pharmaceutical Participant 7

Discussion

This study investigated the existing problems that prevent the data standardization of real-world data for clinical research and further recommendation on each of these problems. Qualitative interviews were conducted on five types of stakeholders which

1
2
3 included hospitals, hospital information system vendors, big data companies,
4 pharmaceutical companies, and regulatory institutions. The wide range of stakeholders
5 were meant to better gauge industry views on the usage of real-world data for clinical
6 research and their thoughts on data standardization.
7

8
9 Difficulty integrating standards to routine clinical process arise from the conflicts
10 between clinical research and routine care source data collection. Data collected
11 routinely is not meant to be expressed without variations or based on a research protocol
12 for data collection; This raises an important question: what is the proper context for the
13 implementation of data standards? The recommendations which include having a better
14 feedback loop for standards and better separation of data standards and documentation
15 may give a signal for the right context of data standardization. If data standards can
16 gradually assist the physician to clarify, summarize, and guide their clinical decision-
17 making process, the chance of data becoming misclassified can be reduced and data
18 quality can be increased. Furthermore, improvements in artificial intelligence and
19 usability of technology can better separate the natural workflow of the two processes
20 by ensuring that standardization can still be done even if there are variations in data
21 input.
22

23
24
25 The problem surrounding the generalizability of research datasets are often the result
26 of differences in the needs for data by the stakeholders. Often, if the data is standardized
27 to fulfill the needs of only a single stakeholder, the data is transformed into a data silo
28 that cannot be reused for other purposes. The recommendation to promote common
29 research data models is to provide a place of cooperation and negotiation between
30 different stakeholders that allows them the room to analyze the data themselves
31 according to their expertise and data access. For retrospective clinical studies,
32 pharmaceutical stakeholders have worked with big data companies utilizing common
33 research data models to conduct feasibility studies without accessing the data; The data
34 quality was usually audited using plausibility measures based on baseline
35 characteristics in previous studies. However, pharmaceutical companies want to
36 eventually increase usage of real-world data as source data for prospective clinical
37 studies which may require more access and better traceability. In studies that have these
38 requirements, recommendations were made to better integrate electronic data capture
39 and real-world data systems to enforce good clinical practice principles that would help
40 alleviate regulatory evaluation and hospital security concerns. The integration would
41 provide for more direct traceability to the source data and eventually lead to the
42 development of interfaces that can audit the data more efficiently.
43

44
45
46 Recently published FDA (Food Drug Administration) guidance on “Data Standards
47 for Drug and Biological Product Submissions containing Real World Data” have also
48 mentioned similar challenges with real world data standardization.²⁰ Notable challenges
49 including differences in source data standards, data exchange formats, and business
50 processes may often result from the difficulty integrating standards with the routine
51 clinical process leading to variations in standards usage. Similarly, research data sets
52 do not readily exist and may go through complicated data transformation processes that
53 can reduce transparency. FDA further mentions real world data sources may have
54 inconsistent data source formats and wide range of methods and algorithms used to
55
56
57
58
59
60

1
2
3 create these datasets that can be hard to audit. To increase auditing capability, FDA
4 suggests that sponsors should try to conform to FDA-supported data standards, such as
5 CDISC, and document the process of data transformation of source data to these
6 standards which may include audit trail, data curation, differences in data definition,
7 and quality control processes. Our study has also found that stakeholders in China are
8 moving towards a similar direction in leveraging more broadly supported data models
9 and working on creating more transparency and traceability in the data transformation
10 process.
11
12
13

14 15 **Conclusion**

16 The qualitative study investigated the barriers that prevent real world data
17 standardization for clinical research based on constructivist grounded theory. A wide
18 range of stakeholders involved in the source data, standardization process, and validity
19 of data for clinical research were interviewed. The main causes were found to be
20 difficulty integrating standards in routine clinical process, lack of generalizable
21 research datasets, and difficulty auditing the standardization process. The main
22 suggestions for these barriers were found to be better feedback loop for standards, better
23 separation of standards and documentation, promoting generalizable clinical research
24 data models, and improving traceability to source data. The key findings in the study
25 make a clear point that data standardization relies on the consensus of many working
26 groups. Data standards cannot be integrated at the source data unless it is customized
27 for practical usage and can only be secondarily applied to source data adequately when
28 the meaning of the source data can be clearly understood.
29
30
31
32
33

34 **Acknowledgements** We thank all individuals who took the time to participate in our
35 interviews.
36

37 **Contributors** JL, CY, and CL designed the study. JL and XL collected the data. CY
38 and JL contacted the respondents. JL, XL, FJ, and WB analyzed the data. JL and XL
39 wrote the first draft of the manuscript. FJ, CY, and CL revised the manuscript. All
40 authors contributed to the interpretation of the data and editing of the manuscript and
41 approved the final manuscript. CY had full access to all data in the study and had final
42 responsibility for the decision to submit for publication.
43
44

45 **Funding Statement** Merck Sharp and Dohme ECT2109014314.

46 **Ethics Approval** Ethical approval was obtained from Peking University Institutional
47 Review board (No. IRB00001052-21081).
48
49

50 **Figure 1 Caption** Barrier and Suggestions in Data Standardization of Real-World Data
51 for Clinical Research
52
53

54 **References**

- 55
56
57 1. Administration USFaD. Use of Real-World Evidence to Support Regulatory
58 Decision-Making for Medical Devices - Guidance for Industry and Food and
59 Drug Administration Staff: Food and Drug Administration, 2017.
60

- 1
2
3 Available: <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/use-real-world-evidence-support-regulatory-decision-making-medical-devices>
4
5
6
7
8 2. Administration USFaD . Framework for FDA’s Real-World Evidence
9 Program: Food and Drug Administration, 2018.
10 Available: <https://www.fda.gov/science-research/science-and-research-special-topics/real-world-evidence>
11
12
13 3. Sun X, Tan J, Tang L, et al. . Real world evidence: experience and lessons
14 from China. *BMJ* 2018;360:j5262. 10.1136/bmj.j5262
15
16 4. Corrigan-Curay J, Sacks L, Woodcock J. Real-World evidence and real-world
17 data for evaluating drug safety and effectiveness. *JAMA* 2018;320:867–8.
18 10.1001/jama.2018.10136
19
20 5. Statistical Information Center of National Health Commission. Conceptual data
21 model of national health and population information [EB / OL] (2020-05-22)
22 [http://www.nhc.gov.cn/wjw/s9497/202006/7e895db195394b36ac4c806f2748c](http://www.nhc.gov.cn/wjw/s9497/202006/7e895db195394b36ac4c806f2748cce3.shtml?tdsourcetag=s_pctim_aiomsg)
23 [ce3.shtml?tdsourcetag=s_pctim_aiomsg](http://www.nhc.gov.cn/wjw/s9497/202006/7e895db195394b36ac4c806f2748cce3.shtml?tdsourcetag=s_pctim_aiomsg)
24
25 6. General Office of the Ministry of Health and the State Administration of
26 Traditional Chinese Medicine. The basic architecture and data standard of
27 electronic medical record (Trial Implementation) [EB / OL] (2009-12-31)
28 <http://www.nhc.gov.cn/cms-search/xxgk/getManuscriptXxgk.htm?id=45414>.
29
30 7. Statistical Information Center of National Health Commission. Specification for
31 electronic medical record sharing document part 1 (2020 version) [EB / OL]
32 (2016-09-29)
33 [http://www.nhc.gov.cn/wjw/s9497/201609/f725ee2635c74ed3a728cd2350953](http://www.nhc.gov.cn/wjw/s9497/201609/f725ee2635c74ed3a728cd2350953bff.shtml)
34 [bff.shtml](http://www.nhc.gov.cn/wjw/s9497/201609/f725ee2635c74ed3a728cd2350953bff.shtml)
35
36 8. General Office of the National Health Commission. Administrative measures
37 for hierarchical evaluation of the application level of electronic medical record
38 system (for Trial Implementation) and evaluation standards (for Trial
39 Implementation) [EB / OL] (2018-12-07)
40 [http://www.nhc.gov.cn/yzygj/s7659/201812/3cae6834a65d48e9bfd783f3c7d5](http://www.nhc.gov.cn/yzygj/s7659/201812/3cae6834a65d48e9bfd783f3c7d54745.shtml)
41 [4745.shtml](http://www.nhc.gov.cn/yzygj/s7659/201812/3cae6834a65d48e9bfd783f3c7d54745.shtml).
42
43 9. Statistical Information Center of the National Health Commission. Standardized
44 maturity evaluation scheme for hospital information interconnection (2020
45 version) [EB / OL] (2020-08-06)
46 [http://www.nhc.gov.cn/mohwsbwstjxxzx/s8553/202008/e80dafa1334c44c38f6](http://www.nhc.gov.cn/mohwsbwstjxxzx/s8553/202008/e80dafa1334c44c38f644602406a4973.shtml)
47 [44602406a4973.shtml](http://www.nhc.gov.cn/mohwsbwstjxxzx/s8553/202008/e80dafa1334c44c38f644602406a4973.shtml).
48
49
50 10. Jin F, Yao C, Yan X, et al. Gap between real-world data and clinical research
51 within hospitals in China: a qualitative study. *BMJ Open*. 2020;10(12):e038375.
52 Published 2020 Dec 29. doi:10.1136/bmjopen-2020-038375
53
54 11. C K. Constructing Grounded Theory: A practical guide through qualitative
55 analysis Kathy Charmaz Constructing Grounded Theory: A practical guide
56 through qualitative analysis Sage 224 £19.99 0761973532 0761973532
57 [Formula: see text]. *Nurse Res* 2006;13:84. 10.7748/nr.13.4.84.s4
58
59 12. Woods P, Gapp R, King MA. Generating or developing grounded theory:
60

- methods to understand health and illness. *Int J Clin Pharm* 2016;38:663–70. 10.1007/s11096-016-0260-2
13. Tong A, Sainsbury P, Craig J. Consolidated criteria for reporting qualitative research (COREQ): a 32-item checklist for interviews and focus groups. *Int J Qual Health Care* 2007;19:349–57. 10.1093/intqhc/mzm042
 14. Setia MS. Methodology series module 5: sampling strategies. *Indian J Dermatol* 2016;61:505–9. 10.4103/0019-5154.190118
 15. Moser A, Korstjens I. Series: practical guidance to qualitative research. Part 3: sampling, data collection and analysis. *Eur J Gen Pract* 2018;24:9–18. 10.1080/13814788.2017.1375091
 16. Whiting LS. Semi-structured interviews: guidance for novice researchers. *Nurs Stand* 2008;22:35–40. 10.7748/ns2008.02.22.23.35.c6420
 17. Peters K, Halcomb E. Interviews in qualitative research. *Nurse Res* 2015; 22:6–7. 10.7748/nr.22.4.6.s2
 18. Britten N. Qualitative interviews in medical research. *BMJ* 1995;311:251–3. 10.1136/bmj.311.6999.251
 19. Rabiee F. Focus-group interview and data analysis. *Proc Nutr Soc* 2004;63:655–60. 10.1079/PNS2004399
 20. Administration USFaD. Use of Real-World Evidence to Support Regulatory Decision-Making for Medical Devices - Guidance for Industry and Food and Drug Administration Staff. Food and Drug Administration, 2017. Available: <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/use-real-world-evidence-support-regulatory-decision-making-medical-devices>

Table 1 Demographics of the participants

Type of Stakeholder (# of Institutions)	Total Number of Participants
Hospital (8)	16
Hospital System Vendor (4)	10
Big Data Company (3)	15
Pharmaceutical (6)	12
Regulatory (4)	9

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

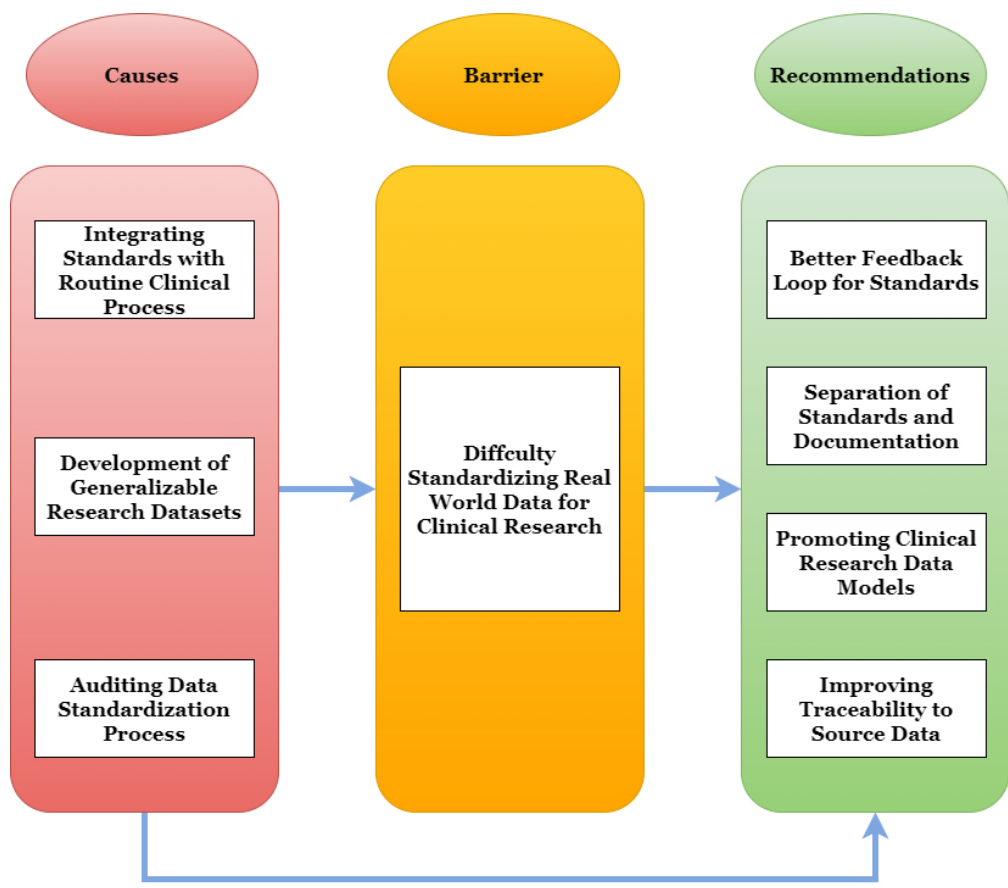


Figure 1 Barrier and Suggestions in Data Standardization of Real-World Data for Clinical Research

289x254mm (72 x 72 DPI)

BMJ Open: first published as 10.1136/bmjopen-2021-059029 on 3 August 2022. Downloaded from <http://bmjopen.bmj.com/> on April 17, 2024 by guest. Protected by copyright.

Consolidated criteria for reporting qualitative studies (COREQ): 32-item checklist

Please indicate in which section each item has been reported in your manuscript. If you do not feel an item applies to your manuscript, please enter N/A.

For further information about the COREQ guidelines, please see Tong *et al.*, 2017:

<https://doi.org/10.1093/intqhc/mzm042>

No.	Item	Description	Section #
Domain 1: Research team and reflexivity			
Personal characteristics			
1.	Interviewer/facilitator	Which author/s conducted the interview or focus group?	
2.	Credentials	What were the researcher's credentials? <i>E.g. PhD, MD</i>	
3.	Occupation	What was their occupation at the time of the study?	
4.	Gender	Was the researcher male or female?	
5.	Experience and training	What experience or training did the researcher have?	
Relationship with participants			
6.	Relationship established	Was a relationship established prior to study commencement?	
7.	Participant knowledge of the interviewer	What did the participants know about the researcher? <i>E.g. Personal goals, reasons for doing the research</i>	
8.	Interviewer characteristics	What characteristics were reported about the interviewer/facilitator? <i>E.g. Bias, assumptions, reasons and interests in the research topic</i>	
Domain 2: Study design			
Theoretical framework			
9.	Methodological orientation and theory	What methodological orientation was stated to underpin the study? <i>E.g. grounded theory, discourse analysis, ethnography, phenomenology, content analysis</i>	
Participant selection			
10.	Sampling	How were participants selected? <i>E.g. purposive, convenience, consecutive, snowball</i>	
11.	Method of approach	How were participants approached? <i>E.g. face-to-face, telephone, mail, email</i>	
12.	Sample size	How many participants were in the study?	
13.	Non-participation	How many people refused to participate or dropped out? What were the reasons for this?	
Setting			
14.	Setting of data collection	Where was the data collected? <i>E.g. home, clinic, workplace</i>	
15.	Presence of non-participants	Was anyone else present besides the participants and researchers?	

16.	Description of sample	What are the important characteristics of the sample? <i>E.g. demographic data, date</i>	
Data collection			
17.	Interview guide	Were questions, prompts, guides provided by the authors? Was it pilot tested?	
18.	Repeat interviews	Were repeat interviews carried out? If yes, how many?	
19.	Audio/visual recording	Did the research use audio or visual recording to collect the data?	
20.	Field notes	Were field notes made during and/or after the interview or focus group?	
21.	Duration	What was the duration of the interviews or focus group?	
22.	Data saturation	Was data saturation discussed?	
23.	Transcripts returned	Were transcripts returned to participants for comment and/or correction?	
Domain 3: analysis and findings			
Data analysis			
24.	Number of data coders	How many data coders coded the data?	
25.	Description of the coding tree	Did authors provide a description of the coding tree?	
26.	Derivation of themes	Were themes identified in advance or derived from the data?	
27.	Software	What software, if applicable, was used to manage the data?	
28.	Participant checking	Did participants provide feedback on the findings?	
Reporting			
29.	Quotations presented	Were participant quotations presented to illustrate the themes / findings? Was each quotation identified? <i>E.g. Participant number</i>	
30.	Data and findings consistent	Was there consistency between the data presented and the findings?	
31.	Clarity of major themes	Were major themes clearly presented in the findings?	
32.	Clarity of minor themes	Is there a description of diverse cases or discussion of minor themes?	

When submitting your manuscript via the online submission form, please upload the completed checklist as a Figure/supplementary file.

If you would like this checklist to be included alongside your article, we ask that you upload the completed checklist to an online repository and include the guideline type, name of the repository, DOI and license in the *Data availability* section of your manuscript.

Developed from: Allison Tong, Peter Sainsbury, Jonathan Craig, Consolidated criteria for reporting qualitative research (COREQ): a 32-item checklist for interviews and focus groups, *International Journal for Quality in Health Care*, Volume 19, Issue 6, December 2007, Pages 349–357, <https://doi.org/10.1093/intqhc/mzm042>

1
2
3
4 List of Institutions:

5 Hospitals:

- 6
7 1. Peking University People's Hospital, Beijing, Beijing, China
8 2. Peking University First Hospital, Beijing, Beijing, China
9 3. First Teaching Hospital of Tianjin University of Traditional Chinese Medicine, Tianjin,
10 Tianjin, China
11 4. Hainan General Hospital, Haikou, Hainan, China
12 5. Boao Evergrande International Hospital, Boao, Hainan, China
13 6. Boao Super Hospital, Boao, Hainan, China
14 7. Boao Yiling Lifecare Center, Hainan, China
15 8. Boao Worldlight Hospital, Hainan, China
16
17

18 Hospital System Vendors:

- 19 1. Haitai International
20 2. Goodwill
21 3. Winning Health
22 4. Orion Health Rhapsody
23
24

25 Big Data Companies:

- 26 1. Yiducloud
27 2. Digital Health China Technologies
28 3. Inspur
29
30

31 Pharmaceutical Companies:

- 32 1. Pfizer
33 2. Tigermed
34 3. AstraZeneca
35 4. Bristol-Meyers Squibb
36 5. Johnson & Johnson
37 6. BeiGene
38
39

40 Regulatory Institutions:

- 41 1. China National Health Development Research Center
42 2. National Medical Products Administration
43 3. China Center for Food and Drug International Exchange
44 4. Hainan Boao Lecheng International Medical Tourism Pilot Zone Administration
45
46
47

48 Questionnaire:

49
50 Hospital:

- 51
52 1. This interview will be recorded and used in a qualitative study, your identity will be
53 concealed to protect your privacy, do we have your full consent in this interview and have
54 your signed information consent form?
55 2. Describe your role and experience facilitating clinical research at the Hospital?
56
57
58
59
60

3. What are the motivating goals of clinical research?
4. How do you determine the research dataset that are needed for your research?
5. Do you think that electronic medical records or routine care data at the hospital are enough to accomplish your research?
6. How do you aggregate and store all data from different hospital systems?
7. How do you implement standards on your data?
8. What areas in your clinical research process do you have to rely on external vendors to help you?
9. How does data sharing happen for medical records inside and outside of the hospital?
10. Beyond clinical research, have you standardized your data for other purposes?

Big Data:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience handling data at the company?
3. Describe your interaction with clients that want to utilize your service for clinical research?
4. How do you organize source data and any standards that you use to do so?
5. How do you manage the variety of standards that are published?
6. How do you transform the data to fit these standards?
7. How do you manage the differences between source data and standards?
8. How do you track the data transformation process?
9. How do you manage the different research projects that need to use real world data?
10. What type of standards are your clients required to fulfill?

Hospital System:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience developing hospital systems? What type of systems does your company produce?
3. What is behind the motivation for hospitals to use more standardized systems?
4. Describe how data is organized and presented for hospital systems and the standards used to create them?
5. How has the usage of standards improved your services?
6. What is the process of negotiating standards usage when customizing your product for clients?
7. How does customization of hospital systems affect standard usage?
8. How do you localize these standards for practical usage?
9. How do you improve communication between your systems or external systems using standards?
10. What areas of hospital systems rely most on existing international standards?

Pharmaceutical:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience in using real world data for clinical research? What types of real-world data do you use as source data for your studies?
3. How do you obtain or access real world data?
4. How does the process of sourcing real world data differ from the traditional data collection for clinical research the most?
5. What standards are used for real world data?
6. What data standards would like to see used for real world data?
7. What standardization methods for real world data are used to produce research data?
8. How do you check whether the data is reliable and what types of data do you think are most reliable?
9. Does real world data meet your research needs?
10. What do you think the standards for real world data necessary to meet evidence requirement by regulatory institutions?

Regulatory:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience in regulating real world studies? What types of real-world data do you see used as source data for these studies?
3. What are the common characteristics of clinical studies using real world data do you often see (study design, phase, purpose)?
4. How does the process of sourcing real world data differ from the traditional data collection for clinical research the most?
5. What are some ethical standards that may be violated in real world data usage for clinical studies?
6. What standards are used for real world data?
7. What data standards are necessary for real world data to meet evidence requirements?
8. What are the main considerations surrounding data standardization?
9. How does real world data meet other regulatory evaluation needs besides clinical trials?
10. How can real world data be better collected to establish a platform in China that can benefit more stakeholders?

BMJ Open

Existing Barriers and Recommendations of Real-World Data Standardization for Clinical Research in China: A Qualitative Study

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-059029.R1
Article Type:	Original research
Date Submitted by the Author:	07-Apr-2022
Complete List of Authors:	Lai, Junkai; Peking University First Hospital, Peking University Clinical Research Institute, Liao, Xiwen; Peking University First Hospital, Peking University Clinical Research Institute Jin, Feifei; Peking University People's Hospital, National Center for Trauma Medicine Wang, Bin; Peking University First Hospital, Peking University Clinical Research Institute Yao, Chen; Peking University First Hospital, Peking University Clinical Research Institute; Hainan Institute of Real World Data Li, Chen; Fourth Military Medical University, Department of Health Statistics; School of Preventive Medicine
Primary Subject Heading:	Health informatics
Secondary Subject Heading:	Qualitative research
Keywords:	QUALITATIVE RESEARCH, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

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

1 **Existing Barriers and Recommendations of Real-World Data Standardization for Clinical**
2 **Research in China: A Qualitative Study**

3
4 Junkai Lai¹, Xiwen Liao¹, Chen Yao^{1,3}, Feifei Jin², Bin Wang¹, Chen Li⁴

5 1 Peking University Clinical Research Institute, Peking University First Hospital, Beijing, China

6 2 National Center for Trauma Medicine, Peking University People's Hospital, Beijing, China

7 3 Hainan Institute of Real World Data, Qionghai, Hainan, China.

8 4 Department of Health Statistics, School of Preventive Medicine, Fourth Military University, Xi'an,
9 Shaanxi, China

10
11 **Correspondence to**

12 Chen Yao

13 Peking University First Hospital, Xicheng District, Beijing 100034, China

14 Tel +86 18610640562

15 Email yaochen@hsc.pku.edu.cn

16
17 **Keywords** real world data; data standards; data standardization; clinical research; qualitative
18 research; China;

19
20 **Word Count**

21 5085

Existing Barriers and Recommendations of Real-World Data Standardization for Clinical Research in China: A Qualitative Study

Abstract

Objective To investigate the existing barriers and recommendations of real-world data (RWD) standardization for clinical research through a qualitative study on different stakeholders.

Design This qualitative study involved five types of stakeholders based on five interview outlines. The data analysis was performed using the constructivist grounded theory analysis process.

Setting 8 Hospitals, 4 Hospital System Vendors, 3 Big Data Companies, 6 Medical Products Companies, and 4 Regulatory Institutions were included.

Participants In total, 62 participants from 25 institutions were interviewed through purposive sampling.

Results The findings showed that the lack of clinical applicability in existing terminology standards, lack of generalizability in existing research databases, and lack of transparency in existing data standardization process were the barriers of data standardization of RWD for clinical research. Expanding coverage of terminology through collecting common terminology, reducing burden in the usage of terminology standards, improving generalizability of RWD for research by using clinical data models, and improving traceability to source data for transparency might be feasible suggestions for solving the current problems.

Conclusions Efficient and reliable data standardization of RWD for clinical research can help generate better evidence used to support regulatory evaluation of medical products. This research suggests expanding coverage of terminology through collecting common terminology, reducing burden in the usage of terminology standards, improving generalizability of RWD for research by using clinical data models, and improving traceability to source data for transparency to guide efforts in data standardization in the future.

Strengths and Limitations of this study:

- Strength: Wide variety of relevant stakeholders on the subject
- Strength: Qualitative understanding of a major industry bottleneck
- Strength: Important recommendations that can guide the direction of the future of the subject
- Limitations: Due to COVID-19, a portion of the interviews were not done in person and may limit the ability to read into the participants response for further exploration of the subject
- Limitations: Recruitment of participants were limited by those that were already exploring the subject and could not be generalized to participants that may not be familiar with the subject. The unselected companies may have different views, which could result in selection bias.

1 Introduction

2 Real world data (RWD) are data relating to patient health status or the delivery of health care
3 collected from a variety of sources such as electronic health records (EHRs)¹⁻⁴. Internationally in
4 the United States (U.S.) and in China, RWD have become increasingly used to support regulatory
5 decision making for drugs and medical devices¹⁻⁵. In September 2019, China's National Medical
6 Products Administration (NMPA) proposed to accelerate the approval process for advanced medical
7 products listed abroad through the collection of RWD from patients using these products in Boao
8 Lecheng Pilot Zone⁶⁻⁷. The proposal has prompted Medical Products companies to conduct clinical
9 research in Boao Lecheng using RWD, specifically electronic medical records (EMR) of patient
10 visits, as evidence for domestic product approval. An example of the first products to leverage the
11 approval process include Johnson & Johnson's femtosecond ophthalmic surgical medical devices
12 which started data collection in October 2019 and was subsequently approved after 6 months⁸. As
13 more products are introduced into Boao Lecheng, there is an imminent need to efficiently translate
14 the data within EMRs to clinical research data.

15
16 The current problem in China is that EMRs constitute a separate system that is not directly connected
17 to electronic data capture (EDC) systems, leading to duplicative and manual transcription of EMR
18 data into the EDC system during clinical research⁹⁻¹⁰. The inefficient process results often in poor
19 data quality due to human error and insufficient source data verification¹¹. Solutions to the issue
20 have been explored by the U.S. Food and Drug Administration (FDA), which includes promoting
21 the direct usage of electronic source data (eSource) within real world data systems such as EHRs
22 for clinical research¹²⁻¹³. In the eSource guidance, interoperability between EHR and EDC systems
23 through the usage of data standards is emphasized. In addition to publishing guidance, initiatives
24 led by the FDA promoted collaboration between standards organizations Health Level Seven (HL7)
25 and Clinical Data Interchange Standards Consortium (CDISC), which have produced solutions that
26 can translate EHR data standards to clinical research data standards¹⁴.

27
28 However, these solutions are not directly translatable to China's context due to differences in the
29 developed data standards and methods used to standardize data. The data standards in China were
30 developed through the Statistical Information Center of the National Health Commission and
31 pushed through government evaluation of hospital information system's meaningful usage¹⁵⁻¹⁸. The
32 first qualitative study on the problem of the gap between RWD and clinical research, found several
33 key domestic problems which included the lack of data standards usage, prevalence of unstructured
34 data, and other data security concerns¹⁹. Similarly, a literature review in China reveals the deterrents
35 of the meaningful usage of RWD for clinical research to include the lack of regulatory
36 implementation of semantic level data standards, unstructured data, and data access²⁰. Therefore,
37 it is important to address the topic of the standardization of RWD for clinical research in China.
38 However, limited literature has addressed the issue and opinions of stakeholders urgently needing
39 to use RWD for clinical research have yet to be collected in China. Therefore, a qualitative study of
40 the relevant stakeholders in China regarding the barriers and suggestions related to the topic is
41 needed.

42 Methods

43 Design

1
2
3 1 Qualitative research allows us to understand a participant's experience through qualitative methods
4 2 of capturing data often through interviews. Grounded theory is a qualitative research method used
5 3 in areas previously unexplored or under explored to inductively generate theory from data grounded
6 4 in the perceptions and concerns of the participant²¹. The method's extensive history in healthcare
7 5 research can be attributed to its systematic process of analysis and stages of coding that allows
8 6 themes to emerge from the data regarding the problems faced by participants and their resolution to
9 7 the problems²². Constructivist grounded theory (CGT) assumes that data are co-constructed through
10 8 the researcher-participant interaction, and the product of analyses is influenced by the interaction of
11 9 the researcher with the data²³⁻²⁴. Studying the underexplored barriers experienced by stakeholders
12 10 and their resolutions in the process of standardizing RWD for clinical research in China is the central
13 11 problem of the study. Therefore, a qualitative research strategy guided by CGT was employed.
14 12

15 13 The research team conducted in-depth interviews with participants. The interviews were conducted
16 14 between September and November 2021. The study is reported following the guidelines of the
17 15 Consolidated Criteria for Reporting Qualitative Research (COREQ) guidelines²⁵.
18 16

17 **Participants Selection**

18 18 The selection of participants was based on their relevance to the type of stakeholders involved in
19 19 the construction of the regional data platform in Boao Lecheng which seeks to efficiently
20 20 standardize RWD for clinical research. The stakeholders include hospitals generating RWD,
21 21 hospital system vendors that install EMRs for hospitals, big data companies that aggregate the
22 22 hospital EMRs onto a data platform, Medical Products companies that consume the data for clinical
23 23 research, and regulatory departments in charge of evaluating the usage of RWD for research. The
24 24 type of stakeholders can be categorized into 3 categories: stakeholders affecting the source data,
25 25 stakeholders affecting the standardization of source data for research, and stakeholders affecting the
26 26 validity of data for regulated clinical research. Hospital and Hospital System Vendors represented
27 27 the first category, Big Data Companies represented the second category, and Medical Products and
28 28 Regulatory Institutions represented the third category.
29 29

30 30 A stratified purposive sampling method was used to select representatives from each of the five
31 31 stakeholder roles²⁶⁻²⁷. Simultaneous data collection and analysis were done to determine when new
32 32 coding information was no longer generated for each role and the interviewing of participants
33 33 stopped²⁸. The resulting number of participants interviewed in the study at information saturation
34 34 included 25 institutions with a total of 62 participants with no dropouts. YC and JL contacted the
35 35 different interviewee and briefed on the subject matter and the objective of investigation before the
36 36 participants agreed to be arranged for an interview. Interviewees represented their own opinions
37 37 based on their experience working at the institution and do not represent the institution. The number
38 38 of participants interviewed by the type of stakeholder is shown in **Table 1**. Detailed list of
39 39 institutions by type of stakeholder is included in the Appendix.
40 40

41 **The inclusion criteria of the interviewees were as follows**

42 **Inclusion criteria**

- 43 44 1. Participants who had extensive experience as a staff member at stakeholder's institution

2. Participants who had experience evaluating RWD for clinical research for the institution

Exclusion criteria

1. Participants who could not sign informed consent form
2. Participants who could not provide at least 45 minutes for an interview

Setting

The research team with training and experience in qualitative methods conducted interviews using a phone or in person. A quiet meeting room was chosen for each interview to allow for better recording of the study data and included only the participant and researchers.

Data Collection

Semi-structured interviews were recorded either over the phone or in person through a phone application with the ability to transcribe audio into text files²⁹⁻³⁰. Field notes were taken to summarize important findings during the interview process, which helped guide later coding. A focus group interview was arranged instead of one-on-one interviews to promote discussion and communication for certain participants³¹. Each interview allowed 60 minutes and basic information including the interview time, place, and interviewee was collected at the beginning. Five sets of interview guides were designed for the five types of stakeholder roles and pilot tested before hand with similar participants that were not included in the study to make the flow of questioning better. Full interview guides were included along with general categories that motivate the questions in the Appendix. The general categories of questions used for each role focus on how the stakeholders affects the data standardization process including at the source, during data standardization to research, and evaluation of research data. The interview questions guided the interviewer in exploring the subject with the participant. Further discussion on the questions or repeat interviews were allowed to explore deeper into the topic or for better clarification. Simultaneous data collection and analysis determined when information saturation had occurred for each role, implying the interviewing of participants ended.

The interviewers were four doctoral students. JL (Male) and XL (Female) were mainly responsible for the interviews, and BW (Male) and FJ (Female) played supportive roles and were mainly responsible for the recordings. The interviewers were trained in a qualitative research course and had experience conducting interviews.

Analysis

All interviews were transcribed to text using the automated transcription software and double checked for each recording by the two interviewers (JL and XL). Coding and memoing were done by three researchers (JL, XL, FJ) who drew on the techniques of constructivist grounded theory while analyzing the data. QSR NVivo V.12 software was used for coding. The team developed a structured coding tree based on the interviews that started with inductive open coding. Once the core categories emerged, deductive selective coding was performed. Memos were used to assist the researchers during the entire analysis process to understand the data, critique the codes, and identify the theoretical categories that the data represented. Open coding was performed independently by two researchers, and the derived core categories were compared in multiple rounds of discussions until all three research members (JL, XL, FJ) agreed. Participants did not provide feedback on the

1
2
3 1 findings.
4 2

5 3 **Patient and public involvement**

6 4 There was no patient or public involvement in this research.
7 5
8 6

9 6 **Results**

10 7 **Barriers and Recommendations in the Standardization of RWD for Clinical Research**

11 8 The CGT framework generated from the three stages of coding and the 62 participants' responses
12 9 were summarized in the flow chart (figure 1). The study found three main barriers and four main
13 10 suggestions. The barriers included lack of clinical applicability in existing terminology standards,
14 11 lack of common data elements in existing databases, and lack of transparency in existing data
15 12 standardization processes. The recommendations included expanding coverage of terminology by
16 13 collecting common terminology, reducing burden in the usage of terminology standards, improving
17 14 applicability of databases using clinical data models, and improving traceability to source data for
18 15 transparency.
19 16

20 17 **Causes**

21 18 *Lack of Clinical Applicability in Existing Terminology Standards*

22 19 The findings showed that hospital and hospital system participants have expressed the lack of
23 20 applicability of terminology standards in the clinical setting. Clinicians expressed that terminology
24 21 standards such as ICD-10 are not granular enough to reflect the diagnosis that they want to make.
25 22 In addition, they expressed that terminology standards often use technical expressions that are not
26 23 commonly used by physicians, making the search process for terminology burdensome. Therefore,
27 24 clinicians expressed that they often use the "other" option to input their own answers. Hospital
28 25 system participants expressed that they often must implement custom made terminology lists
29 26 created by the hospital instead of using default terminology standards to improve the usability of
30 27 the system.
31 28

32 29 "We give our clients default standards to use, but they may feel that the standards do not match their
33 30 needs and will ask us to perform more customizations" – Hospital Information System Vendor
34 31 Participant 1
35 32

36 33 "When implementing standard terminology for the diagnosis field, doctors often just fill in their
37 34 own answers in the "other" option" – Hospital Participant 8
38 35

39 36 *Lack of Common Data Elements in Existing Databases*

40 37 The findings showed that Medical Products companies and regulatory departments expressed that
41 38 the existing RWD databases such as disease specialty databases formed by hospitals are
42 39 standardized to specific research questions and not generalizable to others. Medical Products
43 40 participants expressed that there is substantial variation in the type of available data even when
44 41 standardized. This resulted in the inability to leverage multiple databases together to answer a
45 42 specific clinical research question due to differences in available data and their definitions.
46 43 Regulatory department participants also expressed similar views regarding the applicability of the
47 44 existing RWD databases to support regulatory decision making regarding medical products.
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 1 Currently, the existing data were not organized in a way that could be combined into a generalizable
5 2 research database used to address multiple regulatory questions by different departments.
6 3

7 4 “For feasibility studies, we may look at disease specialty databases. Although data are standardized
8 5 for clinical research, the data elements in these databases are usually very different from each other,
9 6 and we may have to focus on data elements that are more widely available to conduct our studies.”

11 7 -Medical Products Participant 7
12 8

13 9 “Beside our department, other departments are also using RWD in specific datasets. There is
14 10 currently no general platform that can organize RWD to be used by multiple departments to support
15 11 regulatory decision making. Developing such a platform may be in our interest.” – Regulatory
16 12 Participant 3
17 13

14 *Lack of Transparency in Existing Data Standardization Process*

15 15 The findings showed that hospital and Medical Products participants expressed that the data
16 16 standardization process from RWD to clinical research data lacks transparency. Medical Products
17 17 participants expressed that they can use data completeness as well as other metrics to determine the
18 18 quality of the data, but the exact methods used for data standardization are not transparent. In
19 19 addition, they had concerns over the interpretability of standardization methods such as natural
20 20 language processing algorithms in extracting relevant research data and determining whether
21 21 regulatory institutions will accept these methods. Hospital participants also expressed that
22 22 inaccurate data produced by external vendors are difficult to correct or target due to not knowing
23 23 the exact methods used to transform the data. As the producers of research data, big data participants
24 24 expressed that the standardization process requires many steps and teams involved, which can
25 25 reduce its transparency.
26 26

27 27 “The exact methods used for data standardization in producing research databases from RWD are
28 28 not very transparent. My concerns for the usage of hard to interpret artificial intelligence algorithms
29 29 for the extraction and standardization of data are whether regulatory institutions will accept them.”
30 30 – Medical Products Participant 4
31 31

32 32 “When vendors standardize our data into research data, the produced data may sometimes be
33 33 inaccurate. We are not able to understand the methods used in standardization and find the reasons
34 34 why the data may be incorrect.” – Hospital Participant 9
35 35

36 36 “Data standardization may require many teams and communication between many systems, which
37 37 can lead to reduced transparency in the process, making the methods used hard to document
38 38 comprehensively” – Big Data Participant 5
39 39

40 **Suggestions**

41 *Expanding Coverage of Terminology by Collecting Common Terminology*

42 42 The findings showed that big data companies and hospital information system participants
43 43 suggested that the incorporation of their collection of local terminology can improve the coverage
44 44 of the existing terminology standards. Big data participants expressed using RWD to find and
55 56
57 58
59 60

1
2
3
4 1 aggregate colloquial terminology used by clinicians to improve the coverage of terminologies used
5 2 in RWD. Hospital system participants expressed that they have collected practical terminology lists
6 3 from different hospitals, instead of standard terminology lists. In addition, they expressed that these
7 4 local lists are more likely to be used in a clinical setting, because it improved communication and
8 5 could be key to the adoption of terminologies in a clinical setting.
9 6

10 6
11 7 “When working to develop different research databases, our team has incorporated medical experts
12 8 that help us aggregate common terminologies that are synonyms with standard terminology into a
13 9 library. Using the library will help search for relevant RWD.” – Big Data Participant 15
14 10

15 10
16 11 “Standards will get adopted if they can be easily used by our clients. Through our experience
17 12 working with hospitals, we have collected terminology lists that are used, instead of standard
18 13 terminology lists, which improves communication within hospitals.” – Hospital Information System
19 14 Vendor Participant 4
20 15

21 15 22 16 *Reduce Burden in the Usage of Terminology Standards*

23 16
24 17 The findings showed that hospital participants expressed that the efficiency of the usage of data
25 18 standards can be improved by using more automatic methods of terminology standardization.
26 19 Hospital participants expressed various methods used to automatically standardize terminology
27 20 before and after the documentation phase. Before the documentation phase, hospital participants
28 21 suggested that terminology standards can be pre-coordinated with more familiar terminologies
29 22 before usage. After the documentation phase, terminology standards can be post-coordinated
30 23 through natural language processing algorithms that can match local terminologies with standard
31 24 terminology.
32 25

33 25
34 26 “To facilitate the usage of standards during medical documentation, we may recommend more
35 27 familiar terminologies used to display the terminology standards before documentation.” – Hospital
36 28 Participant 9
37 29

38 29
39 30 “Doctors are unfamiliar with the different standards. We will usually work with companies that can
40 31 use better technology such as terminology matching to help us standardize the data after
41 32 documentation.” – Hospital Participant 13
42 33

43 33 44 34 *Improving Applicability of Databases using Clinical Data Models*

45 34
46 35 The findings showed that hospital system and big data participants expressed that the usage of data
47 36 model standards to organize RWD can improve the applicability of RWD to different clinical
48 37 research questions or services. Hospital system participants expressed that the usage of HL7 RIM
49 38 data model can facilitate the efficiency of reusing data for different services including clinical
50 39 decision support services. Big data participants suggested the usage of the OHDSI data model to
51 40 organize their data for the reuse of data to answer different research questions. In addition, they
52 41 suggested that research in different disease areas may require a further extension of the models by
53 42 analyzing where these models fail to capture specific types of data.
54 43

55 43
56 44 “Learning from Huawei’s and Alibaba’s approach to organize their services, we are starting to apply
57 45
58 46
59 47
60 48

1 the HL7 RIM (Health Level 7 Reference Information Model) model to build a middle layer in which
2 our different hospital systems can create their services. Eventually, we would like to use it to support
3 clinical decision support systems” – Hospital Information System Vendor Participant 1

4
5 “When we participate in more clinical studies, we find that the usage of data models such as OHDSI
6 data model can be used to help organize data to answer multiple research questions. However, we
7 may need to extend the data models for more specific diseases by analyzing gap between our schema
8 and the sponsors research case report forms.” – Big Data Participant 5

9 10 *Improving Traceability to Source Data for Transparency*

11 The findings showed that regulatory department and Medical Products participants suggested the
12 improvement in the traceability to source data for better transparency in the data standardization
13 process. Regulatory departments recommended that clinical research involving RWD should adhere
14 to the Good Clinical Practice (GCP) principles which require that research data are traceable to its
15 source data. In addition, aspects of a clinical trial management workflow to authenticate and monitor
16 the quality of the data should be used to increase the confidence in the research data obtained.
17 Medical Products company participants suggested the usage of eSource methods that meet
18 regulatory expectations in terms of auditing the source data to determine the quality of the collected
19 data.

20
21 “The GCP principles should be upheld similarly when using RWD for clinical research. Applying
22 aspects of the clinical trial workflow may be needed to raise the confidence in the quality of RWD
23 collection.” – Regulatory Institution Participant 2

24
25 “We have been searching for eSource tools/companies that can help us collect reliable source data
26 for clinical research that can be easily audited and used as evidence for regulatory approval” -
27 Medical Products Participant 7

28 29 **Discussion**

30 The barriers and recommendations in the standardization of RWD for clinical research is the
31 research question central to the current qualitative study. Through a constructivist grounded theory
32 approach, the study found three main barriers and four main suggestions. The barriers included
33 lack of clinical applicability in existing terminology standards, lack of common data elements in
34 existing databases, and lack of transparency in the existing data standardization process. The
35 recommendations included expanding coverage of terminology by collecting common
36 terminology, reducing burden in the usage of terminology standards, improving applicability of
37 databases using clinical data models, and improving traceability to source data for transparency.
38 The grounded theory used in the paper was applied to address a specific problem regarding the
39 difficulty in RWD standardization for clinical research. The use of the methods in grounded
40 theory were to find the barriers and recommendation to the research problem, with the goal of
41 using the recommendations found to the barriers that similar stakeholders may face in China.

42
43 In this study, the first reason identified was the lack of clinical applicability of current China
44 terminology standards. The current terminology standards do not fit the expressions commonly

1 used by physicians in China and may be burdensome to use. Thus, it is important to promote the
2 collection of common terminology as well as reduce the burden associated with using terminology
3 standards. Internationally, the problem is addressed in many countries through the usage of
4 SNOMED-CT as a comprehensive terminology for clinical application³². The deficiencies of
5 China's EMR standards include its emphasis on the standardization of data elements and limited
6 focus on terminology standards, preventing meaningful exchange of information²⁰. Thus,
7 researchers believed that the localization and implementation of a comprehensive international
8 terminology standard such as SNOMED-CT within EHRs could help represent clinically relevant
9 information comprehensively in China³³. However, previous translation of SNOMED-CT had
10 been insufficient without the collection of terminology synonyms, since physicians did not follow
11 the precise expressions in terminologies³⁴. In contrast, local terminology datasets in China have
12 shown its ability to cover 74.8% of commonly terms used within EHRs³⁵. Therefore, the
13 recommendations to collect local terminology is particularly important to increase the clinical
14 applicability of current terminology standards.

15
16 The other issue regarding clinical applicability of existing terminology standards is the burden
17 associated with its usage. A literature review studying the impact of EHR data structures, such as
18 coding systems, on clinical efficiency found conflicting results with some studies suggesting that
19 structured data made work processes easier while other studies suggesting that coding and
20 entering structured data was slower³⁶. The study further explained that the perceived difficulties
21 might be due to the lack of familiarity with the coding systems. Participants in our study suggested
22 leveraging pre-coordination and post-coordination methods to use terminology standards without
23 depending on a clinician's familiarity with terminology standards. Pre-coordination is a strategy
24 that constrains and maps coding systems to existing local terminology lists, allowing for the usage
25 of local terminology lists without familiarity with external coding systems. A successful
26 implementation of pre-coordination was demonstrated in Hong Kong by binding local
27 terminology, the Hong Kong Clinical Terminology Table (HKCTT), to international terminology
28 standards with the outcome of not influencing regular clinical workflow³⁷. Post-coordination can
29 be applied to existing terminology lists, but here the emphasis is its application to free text by
30 using natural language processing algorithms to extract terms and match them with coding
31 systems. Recent improvements in using NLP showed a 90% accuracy in the extraction and
32 matching of Chinese clinical text terms to SNOMED-CT³⁸. The success of these methods in their
33 respective studies has demonstrated the capability of improving the efficiency of using
34 terminology standards without impacting normal clinical workflow.

35
36 The second reason identified was the lack of generalizability in existing research databases. The
37 lack of generalizability of databases can lead to the limited usage of RWD even after standardization
38 since the databases only address a specific question. Thus, the usage of clinical data models can
39 improve the generalizability of databases by organizing RWD in a consistent and research relevant
40 way to enable the answering of research questions. In the US, the same problem was first discovered
41 in 2008 when met with the technical challenge surrounding the detection of 10 outcomes in 10 drug
42 classes in a network of multiple databases in the Observational Medical Outcomes Partnership
43 (OMOP) research network. The result was the development of a generalizable common data model
44 (CDM) that each database could conform to that would allow for the efficient answering of clinical

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

1 research questions³⁹⁻⁴⁰. In 2021, HL7 and OHDSI (previously OMOP) collectively announced their
2 initiative to create a common data model that integrates data standards common to EHRs with the
3 goal of better organizing EHR data into a clinical research data model⁴¹. Although the usage of
4 common data models in China have not been pushed by the government, the growing usage among
5 big data companies and other research organizations is evident. Confirming the experiences of the
6 participant in the current study, research teams in China have found that even if the same clinical
7 problem is studied, the heterogeneity of cohort studies in terms of variable definition and data
8 collection hinders the integration and sharing of data for clinical research⁴². The problem has been
9 a motivating factor in the review of a suitable international clinical data model that can be used to
10 address the heterogeneity in databases⁴². Application of the OHDSI CDM in China in its first
11 application to study chronic diseases at a single site have now expanded its usage domestically and
12 have also been internationally used to answer COVID-19 treatment questions⁴³⁻⁴⁴. In addition to the
13 application of common data models, translational research and the development of tools to
14 transform related domestic RWD standards, such as HL7 CDA, to common data models, such as
15 OHDSI CDM, are ongoing in Korean and China⁴⁵⁻⁴⁶.

16
17 The final reason is the lack of transparency in the existing data standardization process. The lack
18 of well-documented and understandable methods used in the data standardization process can
19 compromise the reliability of the data for clinical research. Thus, improving traceability of
20 research data to the source data can help evaluate the quality of the standardize data, increase
21 transparency, and meet regulatory expectations. Despite the importance of traceability
22 requirements for regulated clinical research, it remains as a top data standard issue identified by
23 the US FDA in the successful review of submitted data⁴⁷. In response, the US FDA has promoted
24 the use of electronic source data (eSource) including EHRs to enhance the traceability of research
25 data and reduce errors in transcription in several guidance¹²⁻¹³. The implementation of eSource has
26 been researched by the Society of Clinical Data Management to satisfy regulatory expectations
27 regarding data integrity principles⁴⁸. Among the expectations is the emphasis on GCP ALCOA
28 principles including the declaration of source data, usage of standards, real time capture of data,
29 and automatic data quality checks. However, the TransCelerate eSource initiative examined the
30 slow adoption of eSource and found that the main reasons included the lack of standards usage
31 and interoperability between EHRs and EDC systems⁴⁹. In China, researchers have highlighted the
32 need to increase the transparency in the data standardization process by through source data
33 sharing and statistical analysis protocol publishing to increase transparency of standardization
34 methods used⁵⁰. In addition, source data verification, which checks consistency between the data
35 recorded in the database with source data, is promoted with great emphasis by the NMPA where
36 extreme deviations of the source data with research data may lead to legal repercussions⁵¹. To
37 address these issues, suggestions in China were made to develop and utilize an independent
38 eSource platform for the transferring and storing of research source data to guard data integrity
39 and increase transparency. The development and usage of such a platform was tested using real
40 world data collected from the Catalys Precision Laser System medical device real world study in
41 Boao Lecheng and showed great promise in its ability to efficiently transform data while guarding
42 data integrity⁵²⁻⁵³. In 2021, the National Health Commission of China solidified the need for the
43 development and usage of a research source data management platform at medical institutions
44 when they conduct clinical research⁵⁴.

The strength of the study is the selection of a wide and comprehensive range of stakeholder that can better represent the issue in China. Several limitations of this study warrant attention. The participants included specific institutions that were selected to represent the perspective of different stakeholder roles. The unselected companies may have different views, which could result in selection bias. To minimize selection bias, stratified purposive sampling methods were used. Various key institutions were included, and information saturation was assumed to be achieved. In addition, the cultural background and experience of the authors may have influenced the interpretation of the data, although the interviewers had experience and training in conducting qualitative research.

Conclusion

The qualitative study investigated the barriers in RWD standardization for clinical research based on constructivist grounded theory. This study found barriers including lack of clinical applicability in existing terminology standards, lack of common data elements in existing databases, and lack of transparency in existing data standardization process. Expanding coverage of terminology through collecting common terminology, reducing burden in the usage of terminology standards, improving applicability of databases using clinical data models, and improving traceability to source data for transparency may be feasible suggestions for solving the current problems. The findings can be used to promote the development of efficient and reliable methods for the data standardization of RWD for clinical research. Furthermore, the contributions of the study can guide the usage of standards, support the implementation of eSource methods, and facilitate the development of real-world evidence. In the future, we aim to use the suggestions in our study to develop and evaluate eSource tools in China that can standardize RWD for clinical research with efficiency and reliability.

Figure 1 Caption Barrier and Suggestions in Data Standardization of Real-World Data for Clinical Research

Table 1 Demographics of the participants

Type of Stakeholder (# of Institutions)	Total Number of Participants
Hospital (8)	16
Hospital System Vendor (4)	10
Big Data Company (3)	15
Pharmaceutical (6)	12
Regulatory (4)	9

Acknowledgements We thank all individuals who took the time to participate in our interviews.

Contributors JL, CY, and CL designed the study. JL and XL collected the data. CY and JL contacted the respondents. JL, XL, FJ, and WB analyzed the data. JL and XL wrote the first draft of the manuscript. FJ, CY, and CL revised the manuscript. All authors contributed to the interpretation of the data and editing of the manuscript and approved the final manuscript. CY had full access to all data in the study and had final responsibility for the decision to submit for publication.

Funding Statement This work was supported by the Department of Science, Technology, and International Cooperation, National Medical Products Administration.

1 **Competing Interests** None declared.

2 **Patient consent for publication** Not required

3 **Ethics Approval** Ethical approval was obtained from Peking University Institutional Review
4 board (No. IRB00001052-21081).

5 **Data Availability** Data are available upon reasonable request. Study protocol and original data are
6 available on request by emailing the corresponding author.

7

8 **References**

- 9 1 US Food And Drug Administration. Use of Real-World Evidence to Support Regulatory Decision-
10 Making for Medical Devices Guidance for Industry and Food and Drug Administration Staff. In, 2017.
- 11 2 US Food And Drug Administration. REAL-WORLD EVIDENCE PROGRAM. In, 2018.
- 12 3 J Corrigan-Curay, L Sacks, J Woodcock. Real-World Evidence and Real-World Data for Evaluating
13 Drug Safety and Effectiveness. *JAMA* 2018;320(9):867-68.
- 14 4 X Sun, J Tan, L Tang, JJ Guo, X Li. Real world evidence: experience and lessons from China. *BMJ*
15 2018;360:j5262.
- 16 5 National Medical Product Association. Real world evidence supports the guiding principles of drug
17 development and review. In, 2020.
- 18 6 National Development And Reform Commission, National Health Commission, State Administration
19 Of Traditional Chinese Medicine. Implementation measures on supporting the construction of Boao
20 Le Cheng International Medical Tourism Pilot Area. In, 2021.
- 21 7 National Medical Products Association. The State Food and Drug Administration launched the
22 scientific action plan for China's drug supervision. In, 2021.
- 23 8 Johnson Johnson Surgical Vision. A Real-World Evidence Study in China of the Catalys Precision
24 Laser System. In, 2020.
- 25 9 R Blitz, M Dugas. Conceptual Design, Implementation, and Evaluation of Generic and Standard-
26 Compliant Data Transfer into Electronic Health Records. *Appl Clin Inform* 2020;11(3):374-86.
- 27 10 Y Matsumura, A Hattori, S Manabe, et al. Interconnection of electronic medical record with clinical
28 data management system by CDISC ODM. *Stud Health Technol Inform* 2014;205:868-72.
- 29 11 Y Wu, D Yin, K Abbasi. China's medical research revolution. *BMJ* 2018;360:k547.
- 30 12 US Food And Drug Administration. Electronic Source Data in Clinical Investigations. In, 2013.
- 31 13 US Food And Drug Administration. Use of Electronic Health Record Data in Clinical Investigations.
32 In, 2018.
- 33 14 The Office Of The National Coordination For Health Information Technology. Harmonization of
34 Various Common Data Models and Open Standards for Evidence Generation to Support Patient-
35 Centered Outcomes Research. In, 2020.
- 36 15 National Health Commission of the People's Republic of China. Electronic medical record sharing
37 document specification. In, 2020.
- 38 16 National Health Commission of the People's Republic of China. Administrative measures for
39 hierarchical evaluation of application level of electronic medical record system. In, 2018.
- 40 17 National Health Commission of the People's Republic of China. Standardized maturity evaluation
41 scheme for hospital information interconnection. In, 2020.
- 42 18 National Health Committee Of The People'S Republic Of China. Basic architecture and data standard
43 of electronic medical record. In, 2009.

- 1
2
3
4 1 19 F Jin, C Yao, X Yan, et al. Gap between real-world data and clinical research within hospitals in
5 2 China: a qualitative study. *Bmj Open* 2020;10(12):e38375.
- 6 3 20 J Xie, EQ Wu, S Wang, et al. Real-World Data for Healthcare Research in China: Call for Actions.
7 4 *Value Health Reg Issues* 2022;27:72-81.
- 8 5 21 M Byrne. Grounded theory as a qualitative research methodology. *Aorn J* 2001;73(6):1155-56.
- 9 6 22 AL Chapman, M Hadfield, CJ Chapman. Qualitative research in healthcare: an introduction to
10 7 grounded theory using thematic analysis. *J R Coll Physicians Edinb* 2015;45(3):201-05.
- 11 8 23 FK Metelski, J Santos, C Cechinel-Peiter, et al. Constructivist Grounded Theory: characteristics and
12 9 operational aspects for nursing research. *Rev Esc Enferm Usp* 2021;55:e3776.
- 13 10 24 J Mills, A Bonner, K Francis. Adopting a constructivist approach to grounded theory: Implications
14 11 for research design. *Int J Nurs Pract* 2006;12(1):8-13.
- 15 12 25 A Tong, P Sainsbury, J Craig. Consolidated criteria for reporting qualitative research (COREQ): a
16 13 32-item checklist for interviews and focus groups. *Int J Qual Health Care* 2007;19(6):349-57.
- 17 14 26 A Moser, I Korstjens. Series: Practical guidance to qualitative research. Part 3: Sampling, data
18 15 collection and analysis. *Eur J Gen Pract* 2018;24(1):9-18.
- 19 16 27 MS Setia. Methodology Series Module 5: Sampling Strategies. *Indian J Dermatol* 2016;61(5):505-
20 17 09.
- 21 18 28 J Sutton, Z Austin. Qualitative Research: Data Collection, Analysis, and Management. *Can J Hosp*
22 19 *Pharm* 2015;68(3):226-31.
- 23 20 29 K Peters, E Halcomb. Interviews in qualitative research. *Nurse Res* 2015;22(4):6-07.
- 24 21 30 LS Whiting. Semi-structured interviews: guidance for novice researchers. *Nursing standard*
25 22 2008;22(23):35-40.
- 26 23 31 N Britten. Qualitative interviews in medical research. *BMJ* 1995;311(6999):251-53.
- 27 24 32 O Bodenreider, R Cornet, DJ Vreeman. Recent Developments in Clinical Terminologies - SNOMED
28 25 CT, LOINC, and RxNorm. *Yearb Med Inform* 2018;27(1):129-39.
- 29 26 33 Y Zhu, H Pan, L Zhou, et al. Translation and localization of SNOMED CT in China: a pilot study.
30 27 *Artif Intell Med* 2012;54(2):147-49.
- 31 28 34 R Zhang, J Liu, Y Huang, et al. Enriching the international clinical nomenclature with Chinese daily
32 29 used synonyms and concept recognition in physician notes. *BMC Med Inform Decis Mak*
33 30 2017;17(1):54.
- 34 31 35 Y Cheng, T Jiang, L Deng, L Chen, J Ming. Research on the Coverage of Standard Chinese Medical
35 32 Terminology to Practical Application. *Chinese Journal of Health Informatics and Management* 2020.
- 36 33 36 H Forsvik, V Voipio, J Lamminen, et al. Literature Review of Patient Record Structures from the
37 34 Physician's Perspective. *J Med Syst* 2017;41(2):29.
- 38 35 37 K Hung, M Lau, V Fung. Successful Implementation of Terminology Binding in Hong Kong Hospital
39 36 Authority. *Stud Health Technol Inform* 2019;264:1486-87.
- 40 37 38 Y Chen, D Hu, M Li, H Duan, X Lu. Automatic SNOMED CT coding of Chinese clinical terms via
41 38 attention-based semantic matching. *Int J Med Inform* 2022;159:104676.
- 42 39 39 JM Overhage, PB Ryan, CG Reich, AG Hartzema, PE Stang. Validation of a common data model for
43 40 active safety surveillance research. *J Am Med Inform Assoc* 2012;19(1):54-60.
- 44 41 40 PE Stang, PB Ryan, JA Racoosin, et al. Advancing the science for active surveillance: rationale and
45 42 design for the Observational Medical Outcomes Partnership. *Ann Intern Med* 2010;153(9):600-06.
- 46 43 41 Observational Health Data Sciences Informatics. HL7 International and OHDSI Announce
47 44 Collaboration to Provide Single Common Data Model for Sharing Information in Clinical Care and
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- 1 Observational Research. In, 2021.
- 2 42 HX Yue, YL Zhan, F Bian, et al. [Data standard and data sharing in clinical cohort studies]. *Zhonghua*
3 *Liu Xing Bing Xue Za Zhi* 2021;42(7):1299-305.
- 4 43 A Prats-Urbe, AG Sena, L Lai, et al. Use of repurposed and adjuvant drugs in hospital patients with
5 covid-19: multinational network cohort study. *BMJ* 2021;373:n1038.
- 6 44 X Zhang, L Wang, S Miao, et al. Analysis of treatment pathways for three chronic diseases using
7 OMOP CDM. *J Med Syst* 2018;42(12):260.
- 8 45 H Ji, S Kim, S Yi, et al. Converting clinical document architecture documents to the common data
9 model for incorporating health information exchange data in observational health studies: CDA to
10 CDM. *J Biomed Inform* 2020;107:103459.
- 11 46 OMAHA. Mapping with OMOP CDM. In, 2021.
- 12 47 S Hume, S Sarnikar, L Becnel, D Bennett. Visualizing and Validating Metadata Traceability within
13 the CDISC Standards. *AMIA Jt Summits Transl Sci Proc* 2017;2017:158-65.
- 14 48 Society for Clinical Data Management. eSource Implementation in Clinical Research: A Data
15 Management Perspective. In, 2014.
- 16 49 E Kellar, SM Bornstein, A Caban, et al. Optimizing the Use of Electronic Data Sources in Clinical
17 Trials: The Landscape, Part 1. *Ther Innov Regul Sci* 2016;50(6):682-96.
- 18 50 J Xinyao, Z Wenke, Z Junhua, et al. Promote Transparency in Real-World Study. *World Chinese*
19 *Medicine* 2019.
- 20 51 X Yan, C Dong, C Yao. Protecting the accuracy of clinical trial data in China. *BMJ Opinion* 2018.
- 21 52 C Dong, X Yan, R Tian, Z Bian, C YAO. Strengthen the process report of clinical trials, promote full
22 transparency of clinical
23 trials. *Chinese Journal of Evidence-Based Medicine* 2018.
- 24 53 F Jin, C Yao, J Ma, et al. Explore Efficient and Feasible Clinical Real World
25 Data Collection Mode in Hainan Boao Lecheng
26 International Medical Tourism Pilot Zone
27 . *China Food & Drug Administration Magazine* 2021.
- 28 54 National Health Commission of the People's Republic of China. Measures for the administration of
29 clinical research initiated by researchers in medical and health institutions. In, 2021.
- 30

Causes

Barrier

Recommendations

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

Lack of Clinical Applicability in Existing Terminology Standards

Lack of Common Data Elements in Existing Databases

Lack of Transparency in Existing Data Standardization Processes

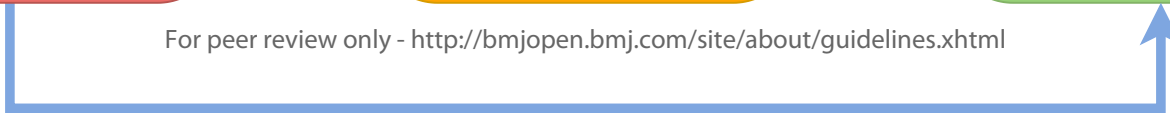
Difficulty Standardizing Real World Data for Clinical Research

Expanding Coverage of Terminology by Collecting Common Terminology

Reducing Burden in the Usage of Terminology Standards

Improving Applicability of Databases using Clinical Data Models

Improving Traceability to Source Data for Transparency



1
2
3
4 List of Institutions:

5 Hospitals:

- 6
7 1. Peking University People's Hospital, Beijing, Beijing, China
8 2. Peking University First Hospital, Beijing, Beijing, China
9 3. First Teaching Hospital of Tianjin University of Traditional Chinese Medicine, Tianjin,
10 Tianjin, China
11 4. Hainan General Hospital, Haikou, Hainan, China
12 5. Boao Evergrande International Hospital, Boao, Hainan, China
13 6. Boao Super Hospital, Boao, Hainan, China
14 7. Boao Yiling Lifecare Center, Hainan, China
15 8. Boao Worldlight Hospital, Hainan, China
16
17

18 Hospital System Vendors:

- 19 1. Haitai International
20 2. Goodwill
21 3. Winning Health
22 4. Orion Health Rhapsody
23
24

25 Big Data Companies:

- 26 1. Yiducloud
27 2. Digital Health China Technologies
28 3. Inspur
29
30

31 Pharmaceutical Companies:

- 32 1. Pfizer
33 2. Tigermed
34 3. AstraZeneca
35 4. Bristol-Meyers Squibb
36 5. Johnson & Johnson
37 6. BeiGene
38
39

40 Regulatory Institutions:

- 41 1. China National Health Development Research Center
42 2. National Medical Products Administration
43 3. China Center for Food and Drug International Exchange
44 4. Hainan Boao Lecheng International Medical Tourism Pilot Zone Administration
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Interview Guide:

Hospital:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience facilitating clinical research at the Hospital?
3. What are the motivating goals of clinical research?
4. How do you determine the data needed for your research? What are the barriers and recommendations?
5. Do you think that electronic medical records or routine care data at the hospital are enough to accomplish your research? What are the barriers and recommendations?
6. How do you use data standards to aggregate and store all data from different hospital systems? What are the barriers and recommendations?
7. What data standards are used and how do you implement data standards during routine data collection? What are the barriers and recommendations?
8. What areas in your clinical research process do you have to rely on external vendors to help you standardize the data and how have you evaluated their data standardization process? What are the barriers and recommendations?
9. How are data standards used to share medical records inside and outside of the hospital? What are the barriers and recommendations?
10. Beyond clinical research, have you standardized your data for other purposes?

Big Data:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience handling data at the company?
3. Describe your interaction with clients that want to utilize your service for clinical research?
4. What type of standards are your clients required to fulfill?
5. How do you use data standards to organize and aggregate source data? What are the barriers and recommendations?
6. How do you use data standards when you transform source data into research datasets? What are the barriers and recommendations?
7. What methods are used to standard source data for clinical research? What are the barriers and recommendations?
8. How do you track and evaluate the quality of the data transformation process? What are the barriers and recommendations?
9. How do you manage the variety of standards that are published? What are the barriers and recommendations?
10. How do you manage the different research projects that need to use real world data? What are the barriers and recommendations?

Hospital System Vendor:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience developing hospital systems? What type of systems does your company produce?
3. What is behind the motivation for hospitals to use data standards?
4. Describe the data standards that are used for hospital systems?
5. How are data standards implemented and customized for the hospital? What are the barriers and recommendations?
6. How do you use data standards to improve internal and external communication at the hospital? What are the barriers and recommendations?

Medical Products Company:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience in using real world data for clinical research? What types of real-world data do you use as source data for your studies?
3. How do you obtain or access real world data?
4. How does the process of sourcing real world data differ from the traditional data collection for clinical research the most?
5. What standards are used for real world data?
6. What data standards would you like to see used for real world data?
7. What standardization methods for real world data are used to produce research data? What are the barriers and recommendations?
8. How do you check whether the data is reliable and what types of data do you think are most reliable? What are the barriers and recommendations?
9. Does real world data meet your research needs? What are the barriers and recommendations?
10. Given regulatory consideration for the usage of real-world data for clinical research, what do you see as the problems in the current methods used to standardize data? What are some recommendations?

Regulatory:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience in regulating real world studies? What types of real-world data do you see used as source data for these studies?
3. What are the common characteristics of clinical studies using real world data do you often see (study design, phase, purpose)?
4. How is real world data used to support regulatory decision making?
5. How does the process of sourcing real world data differ from the traditional data collection for clinical research the most?

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
6. What standards are used for existing real-world data?
 7. What data standards are necessary for the submission of real-world data used to gain product approval?
 8. How can the standardization of real-world data for clinical research meet regulatory expectations? What are the barriers and recommendations?
 9. How can we establish a real-world data platform that can be used for clinical research that can benefit the most stakeholders in China? What are the barriers and recommendations?

For peer review only

General Categories for Interview Questions:

General Category	Hospital	Hospital System Vendor	Big Data Company	Medical Products Company	Regulatory Department
Privacy and Information Consent Statement	1	1	1	1	1
Experience and aim when using RWD for clinical research	2,3,10	2,3	2,3	2,3,4	2,3,4,5
Relevant RWD Standards	7	4	4,9	5,6	6,7
RWD relevance for clinical research	4,5		10	8,9	9
Standardization of RWD at source	6,7,9	5,6	5		
Standardization of RWD for clinical research	8		6,7	7	

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

Reliability of RWD Standardization for clinical research	8		8	8,10	8
--	---	--	---	------	---

For peer review only

BMJ Open: first published as 10.1136/bmjopen-2021-059029 on 3 August 2022. Downloaded from <http://bmjopen.bmj.com/> on April 17, 2024 by guest. Protected by copyright.

Consolidated criteria for reporting qualitative studies (COREQ): 32-item checklist

Please indicate in which section each item has been reported in your manuscript. If you do not feel an item applies to your manuscript, please enter N/A.

For further information about the COREQ guidelines, please see Tong *et al.*, 2017:

<https://doi.org/10.1093/intqhc/mzm042>

No.	Item	Description	Section #
Domain 1: Research team and reflexivity			
Personal characteristics			
1.	Interviewer/facilitator	Which author/s conducted the interview or focus group?	
2.	Credentials	What were the researcher's credentials? <i>E.g. PhD, MD</i>	
3.	Occupation	What was their occupation at the time of the study?	
4.	Gender	Was the researcher male or female?	
5.	Experience and training	What experience or training did the researcher have?	
Relationship with participants			
6.	Relationship established	Was a relationship established prior to study commencement?	
7.	Participant knowledge of the interviewer	What did the participants know about the researcher? <i>E.g. Personal goals, reasons for doing the research</i>	
8.	Interviewer characteristics	What characteristics were reported about the interviewer/facilitator? <i>E.g. Bias, assumptions, reasons and interests in the research topic</i>	
Domain 2: Study design			
Theoretical framework			
9.	Methodological orientation and theory	What methodological orientation was stated to underpin the study? <i>E.g. grounded theory, discourse analysis, ethnography, phenomenology, content analysis</i>	
Participant selection			
10.	Sampling	How were participants selected? <i>E.g. purposive, convenience, consecutive, snowball</i>	
11.	Method of approach	How were participants approached? <i>E.g. face-to-face, telephone, mail, email</i>	
12.	Sample size	How many participants were in the study?	
13.	Non-participation	How many people refused to participate or dropped out? What were the reasons for this?	
Setting			
14.	Setting of data collection	Where was the data collected? <i>E.g. home, clinic, workplace</i>	
15.	Presence of non-participants	Was anyone else present besides the participants and researchers?	

16.	Description of sample	What are the important characteristics of the sample? <i>E.g. demographic data, date</i>	
Data collection			
17.	Interview guide	Were questions, prompts, guides provided by the authors? Was it pilot tested?	
18.	Repeat interviews	Were repeat interviews carried out? If yes, how many?	
19.	Audio/visual recording	Did the research use audio or visual recording to collect the data?	
20.	Field notes	Were field notes made during and/or after the interview or focus group?	
21.	Duration	What was the duration of the interviews or focus group?	
22.	Data saturation	Was data saturation discussed?	
23.	Transcripts returned	Were transcripts returned to participants for comment and/or correction?	
Domain 3: analysis and findings			
Data analysis			
24.	Number of data coders	How many data coders coded the data?	
25.	Description of the coding tree	Did authors provide a description of the coding tree?	
26.	Derivation of themes	Were themes identified in advance or derived from the data?	
27.	Software	What software, if applicable, was used to manage the data?	
28.	Participant checking	Did participants provide feedback on the findings?	
Reporting			
29.	Quotations presented	Were participant quotations presented to illustrate the themes / findings? Was each quotation identified? <i>E.g. Participant number</i>	
30.	Data and findings consistent	Was there consistency between the data presented and the findings?	
31.	Clarity of major themes	Were major themes clearly presented in the findings?	
32.	Clarity of minor themes	Is there a description of diverse cases or discussion of minor themes?	

When submitting your manuscript via the online submission form, please upload the completed checklist as a Figure/supplementary file.

If you would like this checklist to be included alongside your article, we ask that you upload the completed checklist to an online repository and include the guideline type, name of the repository, DOI and license in the *Data availability* section of your manuscript.

Developed from: Allison Tong, Peter Sainsbury, Jonathan Craig, Consolidated criteria for reporting qualitative research (COREQ): a 32-item checklist for interviews and focus groups, International Journal for Quality in Health Care, Volume 19, Issue 6, December 2007, Pages 349–357, <https://doi.org/10.1093/intqhc/mzm042>

BMJ Open

Existing Barriers and Recommendations of Real-World Data Standardization for Clinical Research in China: A Qualitative Study

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-059029.R2
Article Type:	Original research
Date Submitted by the Author:	04-Jul-2022
Complete List of Authors:	Lai, Junkai; Peking University First Hospital, Peking University Clinical Research Institute, Liao, Xiwen; Peking University First Hospital, Peking University Clinical Research Institute Yao, Chen; Peking University First Hospital, Peking University Clinical Research Institute; Hainan Institute of Real World Data Jin, Feifei; Peking University People's Hospital, National Center for Trauma Medicine Wang, Bin; Peking University First Hospital, Peking University Clinical Research Institute Li, Chen; Fourth Military Medical University, Department of Health Statistics; School of Preventive Medicine Zhang, Jun; MSD China Ltd, CORE Liu, Larry; Merck & Co Inc; Weill Cornell Medical College
Primary Subject Heading:	Health informatics
Secondary Subject Heading:	Qualitative research
Keywords:	QUALITATIVE RESEARCH, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

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

1 **Existing Barriers and Recommendations of Real-World Data Standardization for Clinical**
2 **Research in China: A Qualitative Study**

3
4 Junkai Lai¹, Xiwen Liao¹, Chen Yao^{1,3}, Feifei Jin², Bin Wang¹, Chen Li⁴, Jun Zhang⁵, Larry Liu^{6,7}

5 1 Peking University Clinical Research Institute, Peking University First Hospital, Beijing, China

6 2 National Center for Trauma Medicine, Peking University People's Hospital, Beijing, China

7 3 Hainan Institute of Real World Data, Qionghai, Hainan, China.

8 4 Department of Health Statistics, School of Preventive Medicine, Fourth Military University, Xi'an,

9 Shaanxi, China

10 5 MSD R&D (China) Co., Ltd., Beijing, China

11 6 Merck & Co., Inc., Rahway, NJ, USA

12 7 Weill Cornell Medical College, New York, NY, USA

13
14 **Correspondence to**

15 Chen Yao

16 Peking University First Hospital, Xicheng District, Beijing 100034, China

17 Tel +86 18610640562

18 Email yaochen@hsc.pku.edu.cn

19
20 **Keywords** real world data; data standards; data standardization; clinical research; qualitative
21 research; China;

22
23 **Word Count**

24 5238

Existing Barriers and Recommendations of Real-World Data Standardization for Clinical Research in China: A Qualitative Study

Abstract

Objective To investigate the existing barriers and recommendations of real-world data (RWD) standardization for clinical research through a qualitative study on different stakeholders.

Design This qualitative study involved five types of stakeholders based on five interview outlines. The data analysis was performed using the constructivist grounded theory analysis process.

Setting 8 Hospitals, 4 Hospital System Vendors, 3 Big Data Companies, 6 Medical Products Companies, and 4 Regulatory Institutions were included.

Participants In total, 62 participants from 25 institutions were interviewed through purposive sampling.

Results The findings showed that the lack of clinical applicability in existing terminology standards, lack of generalizability in existing research databases, and lack of transparency in existing data standardization process were the barriers of data standardization of RWD for clinical research. Enhancing terminology standards by incorporating locally used clinical terminology, reducing burden in the usage of terminology standards, improving generalizability of RWD for research by using clinical data models, and improving traceability to source data for transparency might be feasible suggestions for solving the current problems.

Conclusions Efficient and reliable data standardization of RWD for clinical research can help generate better evidence used to support regulatory evaluation of medical products. This research suggested enhancing terminology standards by incorporating locally used clinical terminology, reducing burden in the usage of terminology standards, improving generalizability of RWD for research by using clinical data models, and improving traceability to source data for transparency to guide efforts in data standardization in the future.

Strengths and Limitations of this study:

- Strength: Wide variety of relevant stakeholders on the subject
- Strength: Qualitative understanding of a major industry bottleneck
- Strength: Important recommendations that can guide the direction of the future of the subject
- Limitations: Due to COVID-19, a portion of the interviews were not done in person and might limit the ability to read into the participants response for further exploration of the subject
- Limitations: Recruitment of participants were limited to those that were already exploring the subject, which could result in selection bias.

1 Introduction

2 Real world data (RWD) are data relating to patient health status or the delivery of health care
3 collected from a variety of sources such as electronic health records (EHRs) [1-4]. Internationally,
4 especially in the United States (U.S.) and in China, RWD have become increasingly used to support
5 regulatory decision making for drugs and medical devices [1,5]. In September 2019, China's
6 National Medical Products Administration (NMPA) proposed to accelerate the approval process for
7 advanced medical products listed abroad through the collection of RWD from patients using these
8 products in Boao Lecheng Pilot Zone [6-7]. The proposal has prompted Medical Products
9 companies to conduct clinical research in Boao Lecheng using RWD, specifically the patient visit
10 data collected in electronic medical records (EMR), as real-world evidence for domestic product
11 approval. An example of the first products to leverage the approval process included Johnson &
12 Johnson's femtosecond ophthalmic surgical medical devices which started data collection in
13 October 2019 and subsequently gained approval after 14 months [8]. As more products being
14 introduced into Boao Lecheng, there is an imminent need to efficiently translate the data within
15 EMRs to clinical research data.

16
17 A current problem in China is that EMRs constitute a separate system that is not able to be directly
18 connected to electronic data capture (EDC) system used for clinical research data collection, leading
19 to the duplicative and manual transcription of EMR data into the EDC system [9-10]. The inefficient
20 process results in poorer data quality due to the likelihood of human error and insufficient source
21 data verification [11]. Solutions to the issue have been explored by the U.S. Food and Drug
22 Administration (FDA), which includes promoting the direct usage of electronic source data (eSource)
23 from real world data systems for clinical research [12-13]. In the eSource guidance, a key
24 recommendation is to use data standards for the exchange of data to increase interoperability
25 between EHR and EDC systems. In addition, initiatives led by the FDA promoted collaboration
26 between standards organizations like Health Level Seven (HL7) and Clinical Data Interchange
27 Standards Consortium (CDISC), which produced solutions harmonizing the differences between
28 EHR data standards and clinical research data standards [14].

29
30 However, these solutions are not directly translatable to China's clinical research context due to
31 differences in the developed data standards. The data standards in China were developed by the
32 Statistical Information Center of the National Health Commission and used to evaluate the
33 interoperability of hospital information systems [15-18]. The first qualitative study on the problem
34 of the gap between RWD and clinical research found several key problems, which included the lack
35 of data standards usage, prevalence of unstructured data, and data security concerns [19]. Similarly,
36 a literature review in China revealed that meaningful usage of RWD for clinical research is deterred
37 by weak regulatory implementation of semantic level data standards, prevalence of unstructured
38 data, and difficult hospital data access [20]. It is urgently important to address the standardization
39 of RWD for clinical research in China. However, limited literature and stakeholder opinion on the
40 issue exist and have yet to be explored in China. Therefore, our research aimed to explore the
41 barriers and recommendations regarding the standardization of RWD for clinical research in China
42 through a qualitative study conducted on industry-wide stakeholders.

43 Methods

1 **Design**

2 Qualitative research allows us to understand a participant's experience through qualitative methods
3 of capturing data such as the usage of interviews. Grounded theory is a qualitative research method
4 used in research areas that are unexplored or under explored to inductively generate theory from
5 data grounded in the perceptions of the participant [21]. The method's extensive usage in healthcare
6 research can be attributed to its systematic process of coding and analysis that allows important
7 themes to emerge from the data, regarding the problems faced by participants and their resolutions
8 toward these problems [22]. Constructivist grounded theory (CGT) assumes that data are co-
9 constructed through the researcher-participant interaction, and the product of analysis is influenced
10 by the interaction of the researcher with the data [23-24]. This study aimed to examine an
11 underexplored subject, the barriers experienced by stakeholders in the standardization of RWD for
12 clinical research and their recommendations in the context of China. Therefore, a qualitative
13 research strategy guided by CGT was employed.

14
15 The research team conducted in-depth interviews with participants. The interviews were conducted
16 between September and November 2021. The study is reported according to the Consolidated
17 Criteria for Reporting Qualitative Research (COREQ) guidelines [25].

18 **Participants Selection**

19 The selection of participants was based on the type of stakeholders involved in the construction of
20 the regional data platform in Boao Lecheng, which aimed at the standardization of RWD for clinical
21 research. The type of stakeholders included participants from hospitals that generated RWD,
22 hospital system vendors that installed EMRs, big data companies that centralized RWD onto a data
23 platform, medical product companies that accessed RWD for clinical research, and regulatory
24 departments that evaluated the RWD used in clinical research. The type of stakeholders was
25 categorized into 3 general categories: stakeholders that mainly affected the source data, stakeholders
26 that mainly affected the standardization of source data for clinical research, and stakeholders that
27 mainly affected the validity of RWD used for regulated clinical research. Hospital and hospital
28 system vendors represented the first category, big data companies represented the second category,
29 and medical products and regulatory departments represented the third category.

30
31
32 A stratified purposive sampling method was used to select representatives from each of the five
33 stakeholder roles [26-27]. Simultaneous data collection and analysis were conducted to determine
34 when there was no longer new coding information generated for each role and the interviewing of
35 participants stopped [28]. The resulting number of participants interviewed in the study at
36 information saturation included 25 institutions with a total of 62 participants, which included no
37 participant dropouts. YC and JL contacted the interviewees and briefed them on the subject matter
38 of the investigation before the participants agreed to be arranged for an interview. Interviewees
39 represented their own opinions based on their experience working at the institution and do not
40 represent the institution. The number of participants interviewed for each type of stakeholder is
41 shown in **Table 1**. Detailed list of institutions for each type of stakeholder is included in the (See
42 Appendix 1).

43
44 **The inclusion criteria of the interviewees were as follows**

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

Inclusion criteria

1. Participants who had extensive experience as a staff member at stakeholder's institution
2. Participants who had experience evaluating RWD for clinical research for the institution

Exclusion criteria

1. Participants who could not sign informed consent form
2. Participants who could not provide at least 45 minutes for an interview

Setting

The research team with training and experience in qualitative methods conducted interviews using a phone or in person. A quiet meeting room was chosen for each interview to allow for better recording of the study data. Each interview included only the participant and researchers.

Data Collection

Semi-structured interviews were recorded either over the phone or in person through a phone application with the ability to transcribe audio into text files [29-30]. Field notes were taken to summarize important findings during the interview process, which helped guide later coding. A focus group interview was arranged instead of one-on-one interviews to promote discussion and communication for certain participants [31]. Focus groups were used often for hospital and big data teams given the collaborative nature of the work and the tight schedules. Up to three people were involved in a single focus group. Each interview allowed 60 minutes, and basic information, including the interview time, place, and interviewee, was collected at the beginning of the interview. Five sets of interview guides, designed for the five types of stakeholder roles, were pilot tested beforehand with similar participants that were not included in the study to make the flow of questioning better. Full interview guides are included in the appendix along with general categories that motivated these questions (See Appendix 2,3). The general categories of questions used for each role focused on how the stakeholders affected the data standardization process at the source, from the source to research data, and during evaluation at the research data. The interview questions guided the interviewer in exploring the subject with the participant. Further discussion on the questions or repeated interviews were allowed to explore deeper into the topic or for better clarification. Simultaneous data collection and analysis were determined when information saturation had occurred for each role, which implied that the interviewing of participants ended.

The interviewers were four doctoral students. JL (Male) and XL (Female) were mainly responsible for the interviews. BW (Male) and FJ (Female) played supportive roles and were mainly responsible for the recording of interviews. The interviewers were trained in a qualitative research course and had previous experience conducting interviews.

Analysis

All interviews were transcribed to text using the automated transcription software and double checked by the two interviewers (JL and XL). Coding and memoing were done by three researchers (JL, XL, FJ) who drew on the techniques of constructivist grounded theory when they analyzed the data. QSR NVivo V.12 software was used for coding. The team developed a structured coding tree

1 based on the interviews that started with inductive open coding. Once the core categories emerged,
2 deductive selective coding was performed. Memos were used to assist the researchers during the
3 entire analysis process to help them understand the data, critique the codes, and identify the
4 theoretical categories that the data represented. Open coding was performed independently by two
5 researchers, and the derived core categories were compared in multiple rounds of discussions until
6 all three research members (JL, XL, FJ) agreed. Participants did not provide feedback on the
7 findings.

9 **Patient and public involvement**

10 There was no patient or public involvement in this research.

12 **Results**

13 **Barriers and Recommendations in the Standardization of RWD for Clinical Research**

14 The CGT framework generated from the three stages of coding and the 62 participants' responses
15 were summarized in the flow chart (figure 1). The study found three main barriers and four main
16 suggestions. The barriers included lack of clinical applicability in existing terminology standards,
17 lack of common data elements in existing databases, and lack of transparency in existing data
18 standardization processes. The recommendations included enhancing terminology standards by
19 incorporating locally used clinical terminology, reducing burden in the usage of terminology
20 standards, improving applicability of databases using clinical data models, and improving
21 traceability to source data for transparency.

23 **Causes**

24 *Lack of Clinical Applicability in Existing Terminology Standards*

25 The findings showed that hospital and hospital system participants have expressed the lack of
26 applicability of terminology standards in the clinical setting. Clinicians expressed that terminology
27 standards such as ICD-10 are not granular enough to reflect the diagnosis that they want to make.
28 In addition, they expressed that terminology standards often use technical expressions that are not
29 commonly used by physicians, making the search process for terminology burdensome. Therefore,
30 clinicians expressed that they often use the "other" option to input their own answers. Hospital
31 system participants expressed that they often must implement custom made terminology lists
32 created by the hospital instead of using default terminology standards to improve the usability of
33 the system.

34
35 "We give our clients default standards to use, but they may feel that the standards do not match their
36 needs and will ask us to perform more customizations" – Hospital Information System Vendor
37 Participant 1

38
39 "When implementing standard terminology for the diagnosis field, doctors often just fill in their
40 own answers in the "other" option" – Hospital Participant 8

42 *Lack of Common Data Elements in Existing Databases*

43 The findings showed that medical product companies and regulatory departments expressed that the
44 existing RWD databases such as disease specialty databases formed by hospitals are standardized

1 to specific research questions and not generalizable to others. Medical product participants
2 expressed that there is substantial variation in the type of available data even when standardized.
3 This resulted in the inability to leverage multiple databases together to answer a specific clinical
4 research question due to differences in available data and their definitions. Regulatory department
5 participants also expressed similar views regarding the applicability of the existing RWD databases
6 to support regulatory decision making regarding medical products. Currently, the existing data were
7 not organized in a way that could be combined into a generalizable research database used to address
8 multiple regulatory questions by different departments.

9
10 “For feasibility studies, we may look at disease specialty databases. Although data are standardized
11 for clinical research, the data elements in these databases are usually very different from each other,
12 and we may have to focus on data elements that are more widely available to conduct our studies.”
13 -Medical Products Participant 7

14
15 “Beside our department, other departments are also using RWD in specific datasets. There is
16 currently no general platform that can organize RWD to be used by multiple departments to support
17 regulatory decision making. Developing such a platform may be in our interest.” – Regulatory
18 Participant 3

19 20 *Lack of Transparency in Existing Data Standardization Process*

21 The findings showed that hospital and medical product participants expressed that the data
22 standardization process from RWD to clinical research data lacks transparency. Medical product
23 participants expressed that they can use data completeness as well as other metrics to determine the
24 quality of the data, but the exact methods used for data standardization are not transparent. In
25 addition, they had concerns over the interpretability of standardization methods such as natural
26 language processing algorithms in extracting relevant research data and the determination of
27 whether regulatory institutions would accept these methods. Hospital participants also expressed
28 that inaccurate data produced by external vendors are difficult to correct or target due to the
29 unknown methods used to transform the data. As the producers of research data, big data participants
30 expressed that the standardization process requires many steps and teams involved, which can
31 reduce its transparency.

32
33 “The exact methods used for data standardization in producing research databases from RWD are
34 not very transparent. My concerns for the usage of hard to interpret artificial intelligence algorithms
35 for the extraction and standardization of data are whether regulatory institutions will accept them.”
36 – Medical Products Participant 4

37
38 “When vendors standardize our data into research data, the produced data may sometimes be
39 inaccurate. We are not able to understand the methods used in standardization and find the reasons
40 why the data may be incorrect.” – Hospital Participant 9

41
42 “Data standardization may require many teams and communication between many systems, which
43 can lead to reduced transparency in the process and make the methods used hard to document
44 comprehensively” – Big Data Participant 5

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

1 2 **Suggestions**

3 *Enhancing terminology standards by incorporating locally used clinical terminology*

4 The findings showed that big data companies and hospital information system participants
5 suggested that the incorporation of their collection of locally used clinical terminology can improve
6 the coverage of the existing terminology standards. Big data participants expressed the need to add
7 and map RWD terminology found in their databases to standard terminology to enhance current
8 terminology standards. Hospital system participants expressed that they have collected practical
9 terminology lists from different hospitals that are used instead of default standard terminology lists.
10 In addition, they expressed that the choice to use local lists in a clinical setting is to improve better
11 departmental communication and may be a key component in the revision of terminology standards.

12
13 “When working to develop different research databases, our team has incorporated medical experts
14 that help us aggregate common terminologies that are synonyms with standard terminology into a
15 library. Using the library will help search for relevant RWD.” – Big Data Participant 15

16
17 “Standards will get adopted if they can be easily used by our clients. Through our experience
18 working with hospitals, we have collected terminology lists that are used often instead of standard
19 terminology lists due to its ability to improve communication within hospitals.” – Hospital
20 Information System Vendor Participant 4

21 22 *Reduce Burden in the Usage of Terminology Standards*

23 The findings showed that hospital participants expressed that the efficiency of the usage of data
24 standards can be improved by using more automatic methods of terminology standardization.
25 Hospital participants expressed various methods used to automatically standardize terminology
26 before and after the documentation phase. Before the documentation phase, hospital participants
27 suggested that terminology standards can be pre-coordinated with more familiar terminologies
28 before usage. After the documentation phase, terminology standards can be post-coordinated
29 through natural language processing algorithms that can match local terminologies with standard
30 terminology.

31
32 “To facilitate the usage of standards during medical documentation, we may recommend more
33 familiar terminologies used to display the terminology standards before documentation.” – Hospital
34 Participant 9

35
36 “Doctors are unfamiliar with the different standards. We will usually work with companies that can
37 use better technology such as terminology matching to help us standardize the data after
38 documentation.” – Hospital Participant 13

39 40 *Improving Applicability of Databases using Clinical Data Models*

41 The findings showed that hospital system and big data participants expressed that the usage of
42 clinical data model standards to organize RWD can improve the applicability of RWD to different
43 clinical research questions or services. Hospital system participants expressed that the usage of HL7
44 RIM data model can facilitate the reuse of data for different services including clinical decision

1 support services. Big data participants suggested the usage of the OHDSI data model to organize
2 RWD for the answering of different clinical research questions. In addition, they suggested that
3 research in different disease areas may require a further extension of the models by analyzing where
4 these models fail to capture specific types of data.

5
6 “Learning from Huawei’s and Alibaba’s approach to organize their services, we are starting to apply
7 the HL7 RIM (Health Level 7 Reference Information Model) model to build a middle layer in which
8 our different hospital systems can create their services. Eventually, we would like to use it to support
9 clinical decision support systems” – Hospital Information System Vendor Participant 1

10
11 “When we participate in more clinical studies, we find that the usage of data models such as OHDSI
12 data model can be used to help organize data to answer multiple research questions. However, we
13 may need to extend the data models for more specific diseases by analyzing gap between our schema
14 and the sponsors research case report forms.” – Big Data Participant 5

15 16 *Improving Traceability to Source Data for Transparency*

17 The findings showed that regulatory department and medical product participants suggested the
18 improvement in the traceability to source data for better transparency in the data standardization
19 process. Regulatory departments recommended that clinical research involving RWD should adhere
20 to the Good Clinical Practice (GCP) principles which require that research data are traceable to its
21 source data. In addition, aspects of a clinical trial management workflow to authenticate and monitor
22 the quality of the data should be used to increase the confidence in the research data obtained.
23 Medical product company participants suggested the usage of eSource methods that can standardize
24 the transmission of source data and help meet regulatory expectations in terms of auditing the quality
25 of source data used for clinical research.

26
27 “The GCP principles should be upheld similarly when using RWD for clinical research. Applying
28 aspects of the clinical trial workflow may be needed to raise the confidence in the quality of RWD
29 collection.” – Regulatory Institution Participant 2

30
31 “We have been searching for eSource tools/companies that can help us collect reliable source data
32 for clinical research that can be easily audited and used as evidence for regulatory approval” -
33 Medical Products Participant 7

34 35 **Discussion**

36 The barriers and recommendations in the standardization of RWD for clinical research are the
37 research questions central to the current qualitative study. Through a constructivist grounded
38 theory approach, the study found three main barriers and four main suggestions. The barriers
39 included lack of clinical applicability in existing terminology standards, lack of common data
40 elements in existing databases, and lack of transparency in the existing data standardization
41 process. The recommendations included enhancing terminology standards by incorporating locally
42 used clinical terminology, reducing burden in the usage of terminology standards, improving
43 applicability of databases using clinical data models, and improving traceability to source data for
44 transparency. The grounded theory used in the paper was applied to address a specific problem

1
2
3 1 regarding the difficulty in RWD standardization for clinical research. The use of the methods in
4 2 grounded theory was to find the barriers and recommendation to the research problem, with the
5 3 goal of applying the recommendations found to the barriers that similar stakeholders may face in
6 4 China.
7
8
9 5

10 6 In this study, the first reason identified was the lack of clinical applicability of current China
11 7 terminology standards. The current terminology standards do not fit the expressions commonly
12 8 used by physicians in China and may be burdensome to use. Thus, it is important to enhance
13 9 terminology standards by adding locally used clinical terminology as well as reduce the burden
14 10 associated with using terminology standards. Internationally, the problem is addressed in many
15 11 countries through the usage of SNOMED-CT as a comprehensive terminology for clinical
16 12 application [32]. The deficiencies of China's EMR standards include its emphasis on the
17 13 standardization of data elements and limited focus on terminology standards, preventing
18 14 meaningful exchange of information at the semantic level [20]. Thus, researchers believed that the
19 15 localization and implementation of a comprehensive international terminology standard such as
20 16 SNOMED-CT within EHRs could help represent clinically relevant information comprehensively
21 17 in China [33]. However, previous translation of SNOMED-CT had been insufficient without the
22 18 collection of terminology synonyms, since physicians did not follow the precise expressions in
23 19 terminologies [34]. In contrast, local terminology datasets in China showed its ability to cover
24 20 74.8% of commonly terms used within EHRs [35]. Therefore, the recommendations to collect
25 21 local terminology is particularly important to increase the clinical applicability of current
26 22 terminology standards.
27
28
29
30
31
32 23

33 24 The other issue regarding clinical applicability of existing terminology standards is the burden
34 25 associated with its usage. A literature review studying the impact of EHR data structures, such as
35 26 coding systems, on clinical efficiency found conflicting results with some studies suggesting that
36 27 structured data made work processes easier while other studies suggesting that coding and
37 28 entering structured data was slower [36]. The study further explained that the perceived
38 29 difficulties might be due to the lack of familiarity with the coding systems. Participants in our
39 30 study suggested leveraging pre-coordination and post-coordination methods to use terminology
40 31 standards without depending on a clinician's familiarity with terminology standards. Pre-
41 32 coordination is a strategy that constrains and maps coding systems to existing local terminology
42 33 lists, allowing for the usage of local terminology lists without familiarity with external coding
43 34 systems. A successful implementation of pre-coordination was demonstrated in Hong Kong by
44 35 binding local terminology, the Hong Kong Clinical Terminology Table (HKCTT), to international
45 36 terminology standards with the outcome of not influencing regular clinical workflow [37]. Post-
46 37 coordination can be applied to existing terminology lists, but here the emphasis is its application to
47 38 free text by using natural language processing algorithms to extract terms and match them with
48 39 coding systems. Recent improvements in using NLP showed a 90% accuracy in the extraction and
49 40 matching of Chinese clinical text terms to SNOMED-CT [38]. The success of these methods in
50 41 their respective studies has demonstrated the capability of improving the efficiency of using
51 42 terminology standards without impacting normal clinical workflow.
52
53
54
55
56
57
58
59
60

60 44 The second reason identified was the lack of generalizability in existing research databases. The

1 lack of generalizability of databases can lead to the limited usage of RWD even after standardization
2 since the databases only address a specific question. Thus, the usage of clinical data models can
3 improve the generalizability of databases by organizing RWD in a consistent and research relevant
4 way to enable the answering of research questions. In the US, the same problem was first discovered
5 in 2008 when met with the technical challenge surrounding the detection of 10 outcomes in 10 drug
6 classes in a network of multiple databases in the Observational Medical Outcomes Partnership
7 (OMOP) research network. The result was the development of a generalizable common clinical data
8 model (CDM) that each database could conform, allowing for the efficient answering of clinical
9 research questions [39-40]. In 2021, HL7 and OHDSI (previously OMOP) collectively announced
10 their initiative to create a clinical data model that integrated data standards common to EHRs with
11 the goal of better organizing EHR data into a clinical research data model [41]. Although the usage
12 of common data models in China has not been pushed by the government, the growing usage among
13 big data companies and other research organizations is evident. Confirming the experiences of the
14 participant in the current study, research teams in China have found that even if the same clinical
15 problem is studied, the heterogeneity of cohort studies in terms of variable definition and data
16 collection hinders the integration and sharing of data for clinical research [42]. The problem has
17 been a motivating factor in the review of a suitable international clinical data model that can be used
18 to address the heterogeneity in databases [42]. Application of the OHDSI CDM in China in its first
19 application to study chronic diseases at a single site has now expanded to its usage domestically to
20 answer COVID-19 treatment questions using country-wide databases [43-44]. In addition to the
21 application of common data models, translational research and the development of tools to
22 transform related domestic RWD standards, such as HL7 CDA, to common data models, such as
23 OHDSI CDM, are ongoing in Korean and China [45-46].

24
25 The final reason was the lack of transparency in the existing data standardization process. The lack
26 of well-documented and understandable methods used in the data standardization process can
27 compromise the reliability of the data for clinical research. Thus, improving traceability of
28 research data to the source data can help evaluate the quality of the standardize data, increase
29 transparency, and meet regulatory expectations. Despite the importance of traceability
30 requirements for regulated clinical research, it remains as a top data standard issue identified by
31 the US FDA in the successful review of submitted data [47]. In response, the US FDA has
32 promoted the use of electronic source data (eSource) including EHRs to enhance the traceability
33 of research data and reduce errors in transcription in several guidance [12-13]. The
34 implementation of eSource has been researched by the Society of Clinical Data Management to
35 satisfy regulatory expectations regarding data integrity principles [48]. Among the expectations is
36 the emphasis on GCP ALCOA principles including the declaration of source data, usage of
37 standards, real time capture of data, and automatic data quality checks. Further, the TransCelerate
38 eSource initiative examined the slow adoption of eSource and found that the main reasons
39 included the lack of standards usage and interoperability between EHRs and EDC systems [49]. In
40 China, researchers have highlighted the need to increase the transparency of the data
41 standardization process through source data sharing and statistical analysis protocol publishing
42 [50]. In addition, source data verification, which checks consistency between the research data
43 and source data, is promoted with great emphasis by the NMPA, where extreme deviations of the
44 source data with research data may lead to legal repercussions [51]. To address these issues,

1
2
3 1 suggestions in China were made to develop and utilize an independent eSource platform for the
4 2 storage and transmission of research source data to guard data integrity and increase transparency.
5 3 The development and usage of such a platform was tested using real world data collected from the
6 4 Catalys Precision Laser System medical device real world study in Boao Lecheng and showed
7 5 great promise in its ability to efficiently transform data while guarding data integrity [52-53]. In
8 6 2021, the National Health Commission of China solidified the need for the usage of a research
9 7 source data management platform at medical institutions as a requirement for the conduct of
10 8 clinical research [54].
11 9

12 10 The strength of the study was the selection of a wide and comprehensive range of stakeholder that
13 11 better represented the issue in China. Several limitations of this study warranted attention. The
14 12 participants included specific institutions that were selected to represent the perspective of
15 13 different stakeholder roles. The unselected companies may have different views, which could
16 14 result in selection bias. To minimize selection bias, stratified purposive sampling methods were
17 15 used. Various key institutions were included, and information saturation was assumed to be
18 16 achieved. In addition, the cultural background and experience of the authors may have influenced
19 17 the interpretation of the data, although the interviewers had experience and training in conducting
20 18 qualitative research.
21 19

22 20 **Conclusion**

23 21 The qualitative study investigated the barriers in RWD standardization for clinical research
24 22 based on constructivist grounded theory. This study found barriers including lack of clinical
25 23 applicability in existing terminology standards, lack of common data elements in existing databases,
26 24 and lack of transparency in existing data standardization process. Enhancing terminology standards
27 25 by incorporating locally used clinical terminology, reducing burden in the usage of terminology
28 26 standards, improving applicability of databases using clinical data models, and improving
29 27 traceability to source data for transparency may be feasible suggestions for solving the current
30 28 problems. The findings can be used to promote the development of efficient and reliable methods
31 29 for the data standardization of RWD for clinical research. Furthermore, the contributions of the
32 30 study can guide the usage of standards, support the implementation of eSource methods, and
33 31 facilitate the development of real-world evidence. In the future, we aim to use the suggestions in
34 32 our study to develop and evaluate eSource tools in China that can standardize RWD for clinical
35 33 research with efficiency and reliability. Secondly, we aim to use the themes discovered to improve
36 34 communication among relevant stakeholder groups as well as use their collaborative opinion to
37 35 improve the development of data standards that can facilitate the standardization of RWD for
38 36 clinical research.
39 37

40 38 **Figure 1 Caption** Barrier and Suggestions in Data Standardization of Real-World Data for
41 39 Clinical Research
42 40

43 41 **Acknowledgements** We thank all individuals who took the time to participate in our interviews.

44 42 **Contributors** JL, CY, and CL designed the study. JL and XL collected the data. CY and JL
45 43 contacted the respondents. JL, XL, FJ, and WB analyzed the data. JL and XL wrote the first draft
46 44 of the manuscript. FJ, CY, and CL revised the manuscript. JZ contributed to the concept and
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 protocol development, implementation of study, and review and revised the manuscript. LL
 2 contributed to concept development, implementation of study, and review and revised the
 3 manuscript. All authors contributed to the interpretation of the data and editing of the manuscript
 4 and approved the final manuscript. CY had full access to all data in the study and had final
 5 responsibility for the decision to submit for publication.

6 **Funding Statement** This work was supported by Department of Science, Technology, and
 7 International Cooperation, National Medical Products Administration (no award/grant number)
 8 and Merck Sharp & Dohme Corp., a subsidiary of Merck & Co., Inc., Rahway, NJ, USA (no
 9 award/grant number).

10 **Ethics Approval** Ethical approval was obtained from Peking University Institutional Review
 11 board (No. IRB00001052-21081).

12 **Competing Interest Statement** No, there are no competing interests for any author

13
14
15 **Table 1 Demographics of the participants**
16

Type of Stakeholder (# of Institutions)	Total Number of Participants
Hospital (8)	16
Hospital System Vendor (4)	10
Big Data Company (3)	15
Pharmaceutical (6)	12
Regulatory (4)	9

References

- 1 US Food And Drug Administration. Use of Real-World Evidence to Support Regulatory Decision-Making for Medical Devices Guidance for Industry and Food and Drug Administration Staff. In, 2017.
- 2 US Food And Drug Administration. REAL-WORLD EVIDENCE PROGRAM. In, 2018.
- 3 J Corrigan-Curay, L Sacks, J Woodcock. Real-World Evidence and Real-World Data for Evaluating Drug Safety and Effectiveness. *JAMA* 2018;320(9):867-68.
- 4 X Sun, J Tan, L Tang, JJ Guo, X Li. Real world evidence: experience and lessons from China. *BMJ* 2018;360:j5262.
- 5 National Medical Product Association. Real world evidence supports the guiding principles of drug development and review. In, 2020.
- 6 National Development And Reform Commission, National Health Commission, State Administration Of Traditional Chinese Medicine. Implementation measures on supporting the construction of Boao Le Cheng International Medical Tourism Pilot Area. In, 2021.
- 7 National Medical Products Association. The State Food and Drug Administration launched the scientific action plan for China's drug supervision. In, 2021.
- 8 Johnson Johnson Surgical Vision. A Real-World Evidence Study in China of the Catalys Precision Laser System. In, 2020.
- 9 R Blitz, M Dugas. Conceptual Design, Implementation, and Evaluation of Generic and Standard-Compliant Data Transfer into Electronic Health Records. *Appl Clin Inform* 2020;11(3):374-86.
- 10 Y Matsumura, A Hattori, S Manabe, et al. Interconnection of electronic medical record with clinical data management system by CDISC ODM. *Stud Health Technol Inform* 2014;205:868-72.
- 11 Y Wu, D Yin, K Abbasi. China's medical research revolution. *BMJ* 2018;360:k547.
- 12 US Food And Drug Administration. Electronic Source Data in Clinical Investigations. In, 2013.
- 13 US Food And Drug Administration. Use of Electronic Health Record Data in Clinical Investigations. In, 2018.
- 14 The Office Of The National Coordination For Health Information Technology. Harmonization of Various Common Data Models and Open Standards for Evidence Generation to Support Patient-Centered Outcomes Research. In, 2020.
- 15 National Health Commission of the People's Republic of China. Electronic medical record sharing document specification. In, 2020.
- 16 National Health Commission of the People's Republic of China. Administrative measures for hierarchical evaluation of application level of electronic medical record system. In, 2018.
- 17 National Health Commission of the People's Republic of China. Standardized maturity evaluation scheme for hospital information interconnection. In, 2020.
- 18 National Health Committee Of The People'S Republic Of China. Basic architecture and data standard of electronic medical record. In, 2009.
- 19 F Jin, C Yao, X Yan, et al. Gap between real-world data and clinical research within hospitals in China: a qualitative study. *Bmj Open* 2020;10(12):e38375.
- 20 J Xie, EQ Wu, S Wang, et al. Real-World Data for Healthcare Research in China: Call for Actions. *Value Health Reg Issues* 2022;27:72-81.
- 21 M Byrne. Grounded theory as a qualitative research methodology. *Aorn J* 2001;73(6):1155-56.
- 22 AL Chapman, M Hadfield, CJ Chapman. Qualitative research in healthcare: an introduction to

- 1 grounded theory using thematic analysis. *J R Coll Physicians Edinb* 2015;45(3):201-05.
- 2 23 FK Metelski, J Santos, C Cechinel-Peiter, et al. Constructivist Grounded Theory: characteristics and
3 operational aspects for nursing research. *Rev Esc Enferm Usp* 2021;55:e3776.
- 4 24 J Mills, A Bonner, K Francis. Adopting a constructivist approach to grounded theory: Implications
5 for research design. *Int J Nurs Pract* 2006;12(1):8-13.
- 6 25 A Tong, P Sainsbury, J Craig. Consolidated criteria for reporting qualitative research (COREQ): a
7 32-item checklist for interviews and focus groups. *Int J Qual Health Care* 2007;19(6):349-57.
- 8 26 A Moser, I Korstjens. Series: Practical guidance to qualitative research. Part 3: Sampling, data
9 collection and analysis. *Eur J Gen Pract* 2018;24(1):9-18.
- 10 27 MS Setia. Methodology Series Module 5: Sampling Strategies. *Indian J Dermatol* 2016;61(5):505-
11 09.
- 12 28 J Sutton, Z Austin. Qualitative Research: Data Collection, Analysis, and Management. *Can J Hosp
13 Pharm* 2015;68(3):226-31.
- 14 29 K Peters, E Halcomb. Interviews in qualitative research. *Nurse Res* 2015;22(4):6-07.
- 15 30 LS Whiting. Semi-structured interviews: guidance for novice researchers. *Nursing standard*
16 2008;22(23):35-40.
- 17 31 N Britten. Qualitative interviews in medical research. *BMJ* 1995;311(6999):251-53.
- 18 32 O Bodenreider, R Cornet, DJ Vreeman. Recent Developments in Clinical Terminologies - SNOMED
19 CT, LOINC, and RxNorm. *Yearb Med Inform* 2018;27(1):129-39.
- 20 33 Y Zhu, H Pan, L Zhou, et al. Translation and localization of SNOMED CT in China: a pilot study.
21 *Artif Intell Med* 2012;54(2):147-49.
- 22 34 R Zhang, J Liu, Y Huang, et al. Enriching the international clinical nomenclature with Chinese daily
23 used synonyms and concept recognition in physician notes. *BMC Med Inform Decis Mak*
24 2017;17(1):54.
- 25 35 Y Cheng, T Jiang, L Deng, L Chen, J Ming. Research on the Coverage of Standard Chinese Medical
26 Terminology to Practical Application. *Chinese Journal of Health Informatics and Management* 2020.
- 27 36 H Forsvik, V Voipio, J Lamminen, et al. Literature Review of Patient Record Structures from the
28 Physician's Perspective. *J Med Syst* 2017;41(2):29.
- 29 37 K Hung, M Lau, V Fung. Successful Implementation of Terminology Binding in Hong Kong Hospital
30 Authority. *Stud Health Technol Inform* 2019;264:1486-87.
- 31 38 Y Chen, D Hu, M Li, H Duan, X Lu. Automatic SNOMED CT coding of Chinese clinical terms via
32 attention-based semantic matching. *Int J Med Inform* 2022;159:104676.
- 33 39 JM Overhage, PB Ryan, CG Reich, AG Hartzema, PE Stang. Validation of a common data model for
34 active safety surveillance research. *J Am Med Inform Assoc* 2012;19(1):54-60.
- 35 40 PE Stang, PB Ryan, JA Racoosin, et al. Advancing the science for active surveillance: rationale and
36 design for the Observational Medical Outcomes Partnership. *Ann Intern Med* 2010;153(9):600-06.
- 37 41 Observational Health Data Sciences Informatics. HL7 International and OHDSI Announce
38 Collaboration to Provide Single Common Data Model for Sharing Information in Clinical Care and
39 Observational Research. In, 2021.
- 40 42 HX Yue, YL Zhan, F Bian, et al. [Data standard and data sharing in clinical cohort studies]. *Zhonghua
41 Liu Xing Bing Xue Za Zhi* 2021;42(7):1299-305.
- 42 43 A Prats-Uribe, AG Sena, L Lai, et al. Use of repurposed and adjuvant drugs in hospital patients with
43 covid-19: multinational network cohort study. *BMJ* 2021;373:n1038.
- 44 44 X Zhang, L Wang, S Miao, et al. Analysis of treatment pathways for three chronic diseases using

- 1 OMOP CDM. *J Med Syst* 2018;42(12):260.
- 2 45 H Ji, S Kim, S Yi, et al. Converting clinical document architecture documents to the common data
3 model for incorporating health information exchange data in observational health studies: CDA to
4 CDM. *J Biomed Inform* 2020;107:103459.
- 5 46 OMAHA. Mapping with OMOP CDM. In, 2021.
- 6 47 S Hume, S Sarnikar, L Becnel, D Bennett. Visualizing and Validating Metadata Traceability within
7 the CDISC Standards. *AMIA Jt Summits Transl Sci Proc* 2017;2017:158-65.
- 8 48 Society for Clinical Data Management. eSource Implementation in Clinical Research: A Data
9 Management Perspective. In, 2014.
- 10 49 E Kellar, SM Bornstein, A Caban, et al. Optimizing the Use of Electronic Data Sources in Clinical
11 Trials: The Landscape, Part 1. *Ther Innov Regul Sci* 2016;50(6):682-96.
- 12 50 J Xinyao, Z Wenke, Z Junhua, et al. Promote Transparency in Real-World Study. *World Chinese
13 Medicine* 2019.
- 14 51 X Yan, C Dong, C Yao. Protecting the accuracy of clinical trial data in China. *BMJ Opinion* 2018.
- 15 52 C Dong, X Yan, R Tian, Z Bian, C YAO. Strengthen the process report of clinical trials, promote full
16 transparency of clinical
17 trials. *Chinese Journal of Evidence-Based Medicine* 2018.
- 18 53 F Jin, C Yao, J Ma, et al. Explore Efficient and Feasible Clinical Real World
19 Data Collection Mode in Hainan Boao Lecheng
20 International Medical Tourism Pilot Zone
21 . *China Food & Drug Administration Magazine* 2021.
- 22 54 National Health Commission of the People's Republic of China. Measures for the administration of
23 clinical research initiated by researchers in medical and health institutions. In, 2021.
- 24

Causes

Barrier

Recommendations

Lack of Clinical Applicability in Existing Terminology Standards

Lack of Common Data Elements in Existing Databases

Lack of Transparency in Existing Data Standardization Processes

Difficulty Standardizing Real World Data for Clinical Research

Enhancing Term Standards by Incorporating Locally Used Clinical Terms

Reducing Burden in the Usage of Terminology Standards

Improving Applicability of Databases using Clinical Data Models

Improving Traceability to Source Data for Transparency

1
2
3 Appendix 1:

4 List of Institutions:

5 Hospitals:

- 6
7 1. Peking University People's Hospital, Beijing, Beijing, China
8 2. Peking University First Hospital, Beijing, Beijing, China
9 3. First Teaching Hospital of Tianjin University of Traditional Chinese Medicine, Tianjin,
10 Tianjin, China
11 4. Hainan General Hospital, Haikou, Hainan, China
12 5. Boao Evergrande International Hospital, Boao, Hainan, China
13 6. Boao Super Hospital, Boao, Hainan, China
14 7. Boao Yiling Lifecare Center, Hainan, China
15 8. Boao Worldlight Hospital, Hainan, China
16
17

18 Hospital System Vendors:

- 19 1. Haitai International
20 2. Goodwill
21 3. Winning Health
22 4. Orion Health Rhapsody
23
24

25 Big Data Companies:

- 26 1. Yiducloud
27 2. Digital Health China Technologies
28 3. Inspur
29
30

31 Pharmaceutical Companies:

- 32 1. Pfizer
33 2. Tigermed
34 3. AstraZeneca
35 4. Bristol-Meyers Squibb
36 5. Johnson & Johnson
37 6. BeiGene
38
39

40 Regulatory Institutions:

- 41 1. China National Health Development Research Center
42 2. National Medical Products Administration
43 3. China Center for Food and Drug International Exchange
44 4. Hainan Boao Lecheng International Medical Tourism Pilot Zone Administration
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Appendix 2:
Interview Guide:

Hospital:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience facilitating clinical research at the Hospital?
3. What are the motivating goals of clinical research?
4. How do you determine the data needed for your research? What are the barriers and recommendations?
5. Do you think that electronic medical records or routine care data at the hospital are enough to accomplish your research? What are the barriers and recommendations?
6. How do you use data standards to aggregate and store all data from different hospital systems? What are the barriers and recommendations?
7. What data standards are used and how do you implement data standards during routine data collection? What are the barriers and recommendations?
8. What areas in your clinical research process do you have to rely on external vendors to help you standardize the data and how have you evaluated their data standardization process? What are the barriers and recommendations?
9. How are data standards used to share medical records inside and outside of the hospital? What are the barriers and recommendations?
10. Beyond clinical research, have you standardized your data for other purposes?

Big Data:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience handling data at the company?
3. Describe your interaction with clients that want to utilize your service for clinical research?
4. What type of standards are your clients required to fulfill?
5. How do you use data standards to organize and aggregate source data? What are the barriers and recommendations?
6. How do you use data standards when you transform source data into research datasets? What are the barriers and recommendations?
7. What methods are used to standard source data for clinical research? What are the barriers and recommendations?
8. How do you track and evaluate the quality of the data transformation process? What are the barriers and recommendations?
9. How do you manage the variety of standards that are published? What are the barriers and recommendations?
10. How do you manage the different research projects that need to use real world data? What are the barriers and recommendations?

Hospital System Vendor:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience developing hospital systems? What type of systems does your company produce?
3. What is behind the motivation for hospitals to use data standards?
4. Describe the data standards that are used for hospital systems?
5. How are data standards implemented and customized for the hospital? What are the barriers and recommendations?
6. How do you use data standards to improve internal and external communication at the hospital? What are the barriers and recommendations?

Medical Products Company:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience in using real world data for clinical research? What types of real-world data do you use as source data for your studies?
3. How do you obtain or access real world data?
4. How does the process of sourcing real world data differ from the traditional data collection for clinical research the most?
5. What standards are used for real world data?
6. What data standards would you like to see used for real world data?
7. What standardization methods for real world data are used to produce research data? What are the barriers and recommendations?
8. How do you check whether the data is reliable and what types of data do you think are most reliable? What are the barriers and recommendations?
9. Does real world data meet your research needs? What are the barriers and recommendations?
10. Given regulatory consideration for the usage of real-world data for clinical research, what do you see as the problems in the current methods used to standardize data? What are some recommendations?

Regulatory:

1. This interview will be recorded and used in a qualitative study, your identity will be concealed to protect your privacy, do we have your full consent in this interview and have your signed information consent form?
2. Describe your role and experience in regulating real world studies? What types of real-world data do you see used as source data for these studies?
3. What are the common characteristics of clinical studies using real world data do you often see (study design, phase, purpose)?
4. How is real world data used to support regulatory decision making?

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
5. How does the process of sourcing real world data differ from the traditional data collection for clinical research the most?
 6. What standards are used for existing real-world data?
 7. What data standards are necessary for the submission of real-world data used to gain product approval?
 8. How can the standardization of real-world data for clinical research meet regulatory expectations? What are the barriers and recommendations?
 9. How can we establish a real-world data platform that can be used for clinical research that can benefit the most stakeholders in China? What are the barriers and recommendations?

For peer review only

Appendix 3:
General Categories for Interview Questions:

General Category	Hospital	Hospital System Vendor	Big Data Company	Medical Products Company	Regulatory Department
Privacy and Information Consent Statement	1	1	1	1	1
Experience and aim when using RWD for clinical research	2,3,10	2,3	2,3	2,3,4	2,3,4,5
Relevant RWD Standards	7	4	4,9	5,6	6,7
RWD relevance for clinical research	4,5		10	8,9	9
Standardization of RWD at source	6,7,9	5,6	5		
Standardization of RWD for clinical research	8		6,7	7	

Reliability of RWD Standardization for clinical research	8		8	8,10	8
--	---	--	---	------	---

For peer review only

Consolidated criteria for reporting qualitative studies (COREQ): 32-item checklist

Please indicate in which section each item has been reported in your manuscript. If you do not feel an item applies to your manuscript, please enter N/A.

For further information about the COREQ guidelines, please see Tong *et al.*, 2017:

<https://doi.org/10.1093/intqhc/mzm042>

No.	Item	Description	Section #
Domain 1: Research team and reflexivity			
Personal characteristics			
1.	Interviewer/facilitator	Which author/s conducted the interview or focus group?	
2.	Credentials	What were the researcher's credentials? <i>E.g. PhD, MD</i>	
3.	Occupation	What was their occupation at the time of the study?	
4.	Gender	Was the researcher male or female?	
5.	Experience and training	What experience or training did the researcher have?	
Relationship with participants			
6.	Relationship established	Was a relationship established prior to study commencement?	
7.	Participant knowledge of the interviewer	What did the participants know about the researcher? <i>E.g. Personal goals, reasons for doing the research</i>	
8.	Interviewer characteristics	What characteristics were reported about the interviewer/facilitator? <i>E.g. Bias, assumptions, reasons and interests in the research topic</i>	
Domain 2: Study design			
Theoretical framework			
9.	Methodological orientation and theory	What methodological orientation was stated to underpin the study? <i>E.g. grounded theory, discourse analysis, ethnography, phenomenology, content analysis</i>	
Participant selection			
10.	Sampling	How were participants selected? <i>E.g. purposive, convenience, consecutive, snowball</i>	
11.	Method of approach	How were participants approached? <i>E.g. face-to-face, telephone, mail, email</i>	
12.	Sample size	How many participants were in the study?	
13.	Non-participation	How many people refused to participate or dropped out? What were the reasons for this?	
Setting			
14.	Setting of data collection	Where was the data collected? <i>E.g. home, clinic, workplace</i>	
15.	Presence of non-participants	Was anyone else present besides the participants and researchers?	

16.	Description of sample	What are the important characteristics of the sample? <i>E.g. demographic data, date</i>	
Data collection			
17.	Interview guide	Were questions, prompts, guides provided by the authors? Was it pilot tested?	
18.	Repeat interviews	Were repeat interviews carried out? If yes, how many?	
19.	Audio/visual recording	Did the research use audio or visual recording to collect the data?	
20.	Field notes	Were field notes made during and/or after the interview or focus group?	
21.	Duration	What was the duration of the interviews or focus group?	
22.	Data saturation	Was data saturation discussed?	
23.	Transcripts returned	Were transcripts returned to participants for comment and/or correction?	
Domain 3: analysis and findings			
Data analysis			
24.	Number of data coders	How many data coders coded the data?	
25.	Description of the coding tree	Did authors provide a description of the coding tree?	
26.	Derivation of themes	Were themes identified in advance or derived from the data?	
27.	Software	What software, if applicable, was used to manage the data?	
28.	Participant checking	Did participants provide feedback on the findings?	
Reporting			
29.	Quotations presented	Were participant quotations presented to illustrate the themes / findings? Was each quotation identified? <i>E.g. Participant number</i>	
30.	Data and findings consistent	Was there consistency between the data presented and the findings?	
31.	Clarity of major themes	Were major themes clearly presented in the findings?	
32.	Clarity of minor themes	Is there a description of diverse cases or discussion of minor themes?	

When submitting your manuscript via the online submission form, please upload the completed checklist as a Figure/supplementary file.

If you would like this checklist to be included alongside your article, we ask that you upload the completed checklist to an online repository and include the guideline type, name of the repository, DOI and license in the *Data availability* section of your manuscript.

Developed from: Allison Tong, Peter Sainsbury, Jonathan Craig, Consolidated criteria for reporting qualitative research (COREQ): a 32-item checklist for interviews and focus groups, International Journal for Quality in Health Care, Volume 19, Issue 6, December 2007, Pages 349–357, <https://doi.org/10.1093/intqhc/mzm042>