



RWD Guidelines for Programming and Analysis Processes

Contents

1. Introduction	1
1.1. Real-World Data and Evidence	1
1.2. Purpose of Real-World Data	2
1.3. Scope	3
1.4. Types of RWE Studies	3
1.4.1. Observational Studies	3
1.4.2. Trials in Clinical Practice Settings	4
2. Planning and Study Set-Up	4
2.1. Framing of the Research Question	5
2.2. Feasibility Assessment	5
2.3. Engagement with Regulatory Health Authorities	5
2.4. Protocol Development and Review	5
2.5. Real-World Data Checklist	6
2.5.1. Purpose of the Real-World Data Checklist	6
3. RWD Study Council	6
3.1. Purpose	6
3.2. RACI Matrix	6
3.3. Contributors/Roles	7
3.4. General RACI Table	7
3.4.1. RACI Table Examples	8
3.4.2. Conclusion	9
4. Data Provenance, Data Ethics and Data Privacy	9
4.1. Data Provenance	9
4.2. Data Ethics	9
4.3. Data Privacy	10
4.3.1. Adhering to Data Privacy Regulations	10
4.3.2. Statistical Disclosure Methods	10
4.3.3. Data Privacy Methods	11
5. Vendor Engagement	11
5.1. Selection	11
5.2. Implementation	11
5.3. After Implementation	11
5.4. Regulatory Compliance	12
5.5. Case Studies	12
5.5.1. Case Study 1: Vendor Engagement in a Real-World Data Study	12
5.5.2. Case Study 2: Vendor Engagement Process for the FDA-Approved RWE Study	12
6. Fit for Purpose Assessment	13
6.1. Scenario: Designing a Study Using RWD	13
6.2. Defining a Hypothetical Target Trial (HTT)	13
7. Analysis and Submission	14
7.1. Statistical Considerations	14
7.2. Confounding and Biases	14
7.3. Methods to Address Confounding in RWD	15
7.4. Submission of RWD: Current Regulatory Landscape	15
7.5. How the Regulatory Submission Process for RCT Translates into RWD	15
7.6. Submission of RWD: What the Future Holds	16
8. Glossary	17
9. Disclaimer	17
10. Appendices	17
Appendix 1.1: Checklist	17
Appendix 1.2: Case Examples	18
11. References	19
12. Project Contact Information	21
13. Acknowledgements	21

Revision History

Version	Date	Summary
1.0	20 March 2026	Initial version

1. Introduction

As per the FDA,¹¹ real-world data (RWD) is data relating to patient health status and/or the delivery of healthcare routinely collected from a variety of sources. Examples include data derived from electronic health records (EHRs), claims and/or billing data, product and/or disease registry data, and other data sources that inform on health status (e.g. data collected from wearables, patient-generated data).

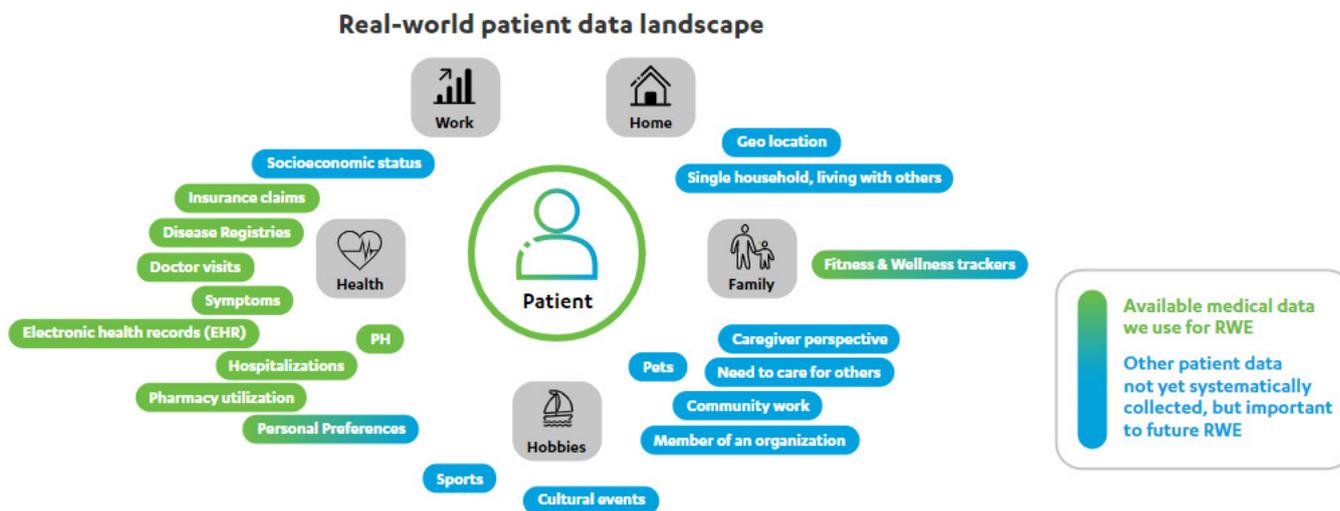


Figure 1.1: Real-world patient data landscape. Source: Janssen RWD Playbook, Janssen R&D [not publicly available]

Real-world evidence (RWE), as defined by the FDA, is clinical evidence regarding the use and potential benefits or risks of a medical product derived from analysis of RWD. Similar definitions have been used by regulatory bodies such as the EMA.

1.1. Real-World Data and Evidence

Randomised controlled trials (RCTs) are considered the gold standard for drug approvals and label claims; however, the demand for RWE is on the rise. Factors contributing to the increasing importance of RWE include:^{1,2}

- New regulatory initiatives
- Access to patient data through medical record databases and disease registries
- Increased interest in:
 - Patient-specific benefits
 - Providing cost-effective large-scale means of monitoring effectiveness and safety (where randomised trials may not be feasible)
- Bridging the evidence gap between clinical research and practice.

RWD has been used successfully in the following cases where RCTs are challenging:

- Rare diseases
- Paediatric studies
- Where randomisation is not possible due to toxicity (oncology)
- Ethical issues with continuing placebo (e.g. Covid-19 vaccine)

Global regulators, such as the US Food & Drug Administration (FDA), Japan's Pharmaceutical and Medical Devices Agency (PMDA), the European Medicines Agency (EMA), the UK's Medicines and Healthcare products Regulatory Agency (MHRA)

and China's National Medical Products Agency (NMPA) are increasingly interested in leveraging the potential of RWD to complement TCTs with RWE to support regulatory decision-making across the product lifecycle.^{1,3} The FDA is accepting observational data to support efficacy determinations and has already issued approvals of new indications for approved drugs.^{1,4} The EMA is assessing the use of registry data^{1,5} and other RWD to support pre- and post-market authorisation activities shown in Figure 1.2. As the healthcare landscape continues to evolve, it is imperative that organisations evolve with it and functions within R&D act proactively to provide the best solutions for supporting RWE activities in line with regulatory expectations. Functions such as Data Management, Programming, and Statistics play a key role in an organisation's ability to leverage and continue to build on cross-functional RWE capabilities and competencies, integrating valuable learnings, new perspectives and best practices across the RWE portfolio.

This paper is a guide for the statistical programming community to achieving these objectives in a pragmatic, efficient and productive manner towards the end goal of serving patients by bringing new treatments to the market.

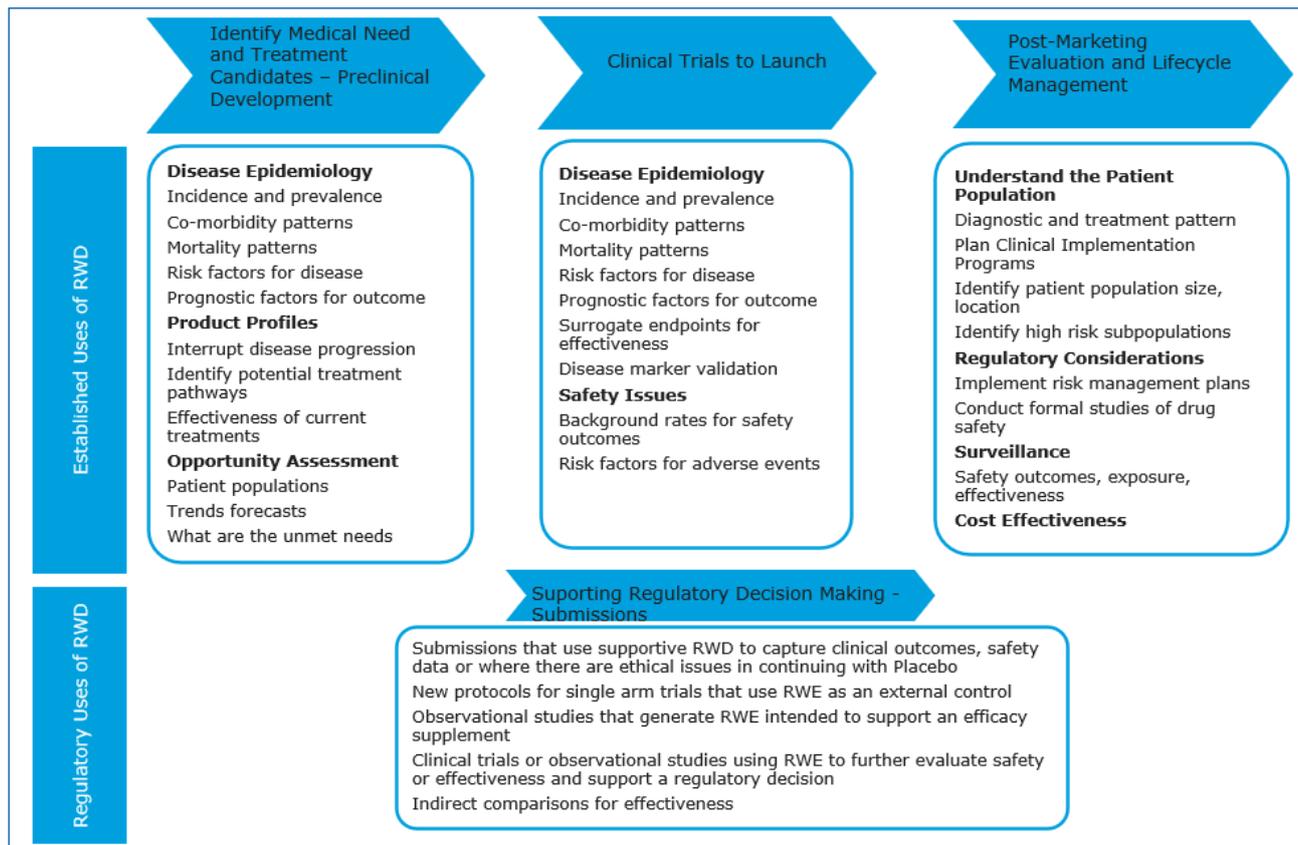


Figure 1.2: Practical uses of RWD during drug development through post-marketing. Source: Janssen RWD Playbook, Janssen R&D [not publicly available]

1.2. Purpose of Real-World Data

Randomised controlled trials (RCTs) are considered the gold standard in clinical research for evaluating the efficacy and safety of interventions. One key strength of RCTs is randomisation, which helps minimise selection bias, ensuring both known and unknown confounding variables are evenly distributed between groups, thereby strengthening the internal validity of the study. By controlling for confounders through randomisation, RCTs provide a strong foundation for making causal inferences about the effects of interventions, which makes them essential for informing clinical practice guidelines, regulatory decisions and healthcare policies. Despite their strengths, RCTs also have limitations, such as high costs, ethical considerations and potential challenges in generalising findings to broader populations.

Here are some of the advantages to integrating RWD with RCTs:

- **Broader patient presentation:** RCTs often have strict eligibility criteria – limiting the diversity and representation of patients – hence may not fully capture how a drug performs in diverse patient groups or under real-world conditions. This can lead to limited applicability of findings. By integrating RWD – including data from real-world clinical practice – a broader range of patients with varying demographics, comorbidities and disease severity can be included in analyses. This enhances the generalisability of findings to real-world populations.

- **Long-term safety and effectiveness:** RCTs are typically conducted over a limited timeframe with a focus on short-term outcomes. Therefore, relying solely on RCTs may not provide sufficient data on rare adverse events or long-term outcomes. RWD, especially from longitudinal studies or post-marketing surveillance, provides valuable insights into the long-term safety and effectiveness of drugs in routine clinical practice. This helps in understanding the drug’s performance over extended periods and in real-world settings.
- **Comprehensive evidence generation:** Integrating RWD with RCTs allows for a more comprehensive evaluation of a drug’s efficacy, safety and real-world impact. RWD can complement RCT data by providing evidence on patient outcomes, treatment patterns, adherence, and healthcare resource use.
- **Cost effectiveness and efficiency:** Leveraging RWD sources such as electronic health records and claim databases alongside RCTs can optimise resource allocation and, in the long term, reduce the cost and time for drug development. RWD can facilitate post-market research and support label expansions by providing supplementary evidence.
- **External controlled arm:** Despite RCTs having obvious appeal in clinical research, they have some well-known limitations. In rare diseases or in diseases where no effective standard-of-care treatments are available, it is not often feasible or ethical to recruit patients to control groups. An uncontrolled, single-arm trial where all participants receive the investigational treatment is more appealing in these scenarios. However,

without an internal control group, assessments are performed by making indirect comparisons, which may be suboptimal.

Using a comparator based on data collected outside of a study – referred to as an external control or synthetic control group – could offer a compromise between uncontrolled trials and RCTs in certain contexts.¹⁶ An external control group could consist of patients treated earlier (sometimes referred to as a historical control) or patients treated during the same period but in a different setting (sometimes referred to as a contemporaneous control).

In summary, combining RWD with RCTs in drug development enhances patient representation, provides insights into long-term outcomes and real-world impact, supports comprehensive evidence generation, and improves cost-effectiveness. This integrated approach can lead to more informed and robust evaluations of drug efficacy and safety, benefiting patients, healthcare providers, and stakeholders involved in drug development and healthcare delivery.

The purposes of using RWD as part of a regulatory submission/decision may include the following:

- To provide evidence in support of the effectiveness or safety of a new product approval
- To provide evidence in support of labelling changes for an approved drug, including:
 - Adding or modifying an indication
 - Change in dose, dose regimen, or route of administration
 - Use in a new population
 - Adding comparative effectiveness information
 - Adding safety information
 - Other labelling change
- To be used as part of a post-marketing requirement to support a regulatory decision:
 - Post-approval safety study
 - Post-approval effectiveness study

1.3. Scope

This white paper is for statistical programmers, especially those who are new to RWD/RWE, who need guidance on recommended best practices and/or minimum requirements for supporting RWE activities that may lead to regulatory submission. Although the primary audience is statistical programmers, other functions typically associated with RCTs may also benefit. In general, 'non-regulatory submission' activities involving RWD are not in scope when it comes to applying the full set of requirements as recommended by regulatory health authorities (RHAs). This may include:

- Exploratory analyses (for internal decision-making/publications)
- Support with analyses for manuscripts and publications
- Market access health technology assessment (HTA) submissions.

This paper chiefly discusses the guidance published by the US FDA for submission of RWD to support marketing of new drug applications (NDAs) or biological products (BLAs). However, the topics discussed are relevant to other regulatory health authorities when it comes to preparing submissions of RWD/RWE for marketing applications.

This paper is organised in a way to mimic the processes that are typical of an RCT submission to RHAs. This way, the reader will have a general understanding of the processes related to RWE studies in comparison to processes relevant to RCTs.

1.4. Types of RWE Studies

The terms RWD and RWE encompass a broad range of study designs, providing more specifics about data sources and the nature of the study. RWD sources (e.g. registries, EHRs, administrative and medical claims databases) can be used for data collection and for analysis infrastructure to support many types of study designs for developing RWE, including randomised trials (e.g. large sample trials, pragmatic clinical trials) and observational studies (prospective or retrospective).

1.4.1. Observational Studies

Observational studies are non-interventional clinical study designs whereby patients receive treatment during routine medical care not guided by research protocols, even if laboratory or imaging procedures are done per protocol.

A non-interventional study is a type of study in which patients receive the marketed drug of interest during routine medical practice and are not assigned to an intervention according to a protocol. Examples of non-interventional study designs include:

- Observational cohort studies, in which patients are identified as belonging to a study group according to the drug/s received or not received during routine medical practice, and subsequent biomedical or health outcomes are identified
- Case control studies, in which patients are identified as belonging to a study group based on having or not having a health-related biomedical or behavioural outcome, and antecedent treatments received are identified.

A **retrospective cohort study** is a study that identifies the population and determines the exposure/treatment from historical data (i.e. data generated before the initiation of the study and therefore after the outcome events have occurred). The variables and outcomes of interest are determined when the study is designed.

This type of study includes database research, review of records, or analysis of electronic healthcare records, where all the events of interest have already happened, and making secondary use of existing data that was not collected for the study purposes.

This type of study can stand alone or be combined with prospective data capture, creating hybrid designs that are more time and cost-efficient than complete de novo data capture.

In a **prospective cohort study**, the population of interest is identified at the start of the study, and patients are enrolled before the occurrence of outcome events and followed prospectively over time. The start of the study is defined as the time at which the research protocol for the study question is initiated. A **registry** is defined as a prospective, non-interventional, organised collection of human data within a particular disease, group or other 'at-risk' special patient

population (e.g. cancer, pregnancy, organ transplant) with design characteristics as follows:

- A systematic collection of defined events or exposures
- Conducted in a defined population in one or more specific geographic areas
- The data collected could be either for a defined period or indefinitely

The questions typically addressed in registries range from purely descriptive questions aimed at understanding the characteristics of patients who develop the disease and how the disease generally progresses, to highly focused questions that support decision-making. Registries focused on determining clinical effectiveness or cost-effectiveness or assessing safety or harm are generally hypothesis-driven and concentrate on evaluating the effects of specific treatments on patient outcomes.

A **registry-based study** is an investigation of a research question using the infrastructure of new or existing registries for patient recruitment and data collection. A registry-based study may be a clinical trial or a non-interventional study and apply primary data collection and/or secondary use of data collected in a patient registry for a purpose other than the given study.

Rare diseases represent a highly heterogeneous group of disorders with high phenotypic and genotypic diversity within individual conditions. Due to the small number of people affected, there are unique challenges in understanding rare diseases and drug development for these conditions, including patient identification and recruitment, trial design, and costs. Natural history data and RWD play significant roles in defining and characterising disease progression, final patient populations, novel biomarkers, genetic relationships, and treatment effects. A natural history study is a preplanned observational study intended to track the course of the disease and identify demographic, genetic, environmental and other variables (e.g. treatment modalities, concomitant medications) that correlate with the disease's development and outcomes. Natural history studies are likely to include patients receiving the current standard of care and/or emergent care, which may alter some manifestations of the disease. Disease registries are a frequent platform to acquire data for natural history studies. For rare diseases, natural history studies play an important role in identifying appropriate patient populations and clinical outcome assessments and biomarkers, and in designing externally controlled studies. Beyond their role in drug development, natural history studies may also benefit patients with rare diseases by establishing communication pathways, identifying disease-specific centres of excellence, facilitating the understanding and evaluation of current standard-of-care practices, evaluating signs and symptoms of a disease to improve diagnosis, and identifying ways to improve patient care. Patients included in natural history studies may be used as historical controls for studies that lack internal control, thus allowing the effectiveness of the study treatment to be determined.

1.4.2. Trials in Clinical Practice Settings

Pragmatic clinical trials (PCTs), sometimes called practical clinical trials, are designed to evaluate the comparative effectiveness of interventions within routine clinical settings. These trials are 'pragmatic' because they focus on understanding how interventions work in real life, as opposed to 'explanatory', where the goal is to determine if and how an intervention works. Key aspects of PCTs are broad population inclusion, study design, and data collection procedures that minimally disrupt routine clinical care encounters, and an emphasis on patient-centred health outcomes. Pragmatic randomised trials (PRTs) represent a hybrid between traditional randomised controlled clinical trials (which are the gold standard for regulatory decision-making) and pragmatic, observational research studies, which are often used to generate real-world evidence. A well-designed PRT that maximises external validity but also controls confounding (including selection bias) to maintain high levels of internal validity could theoretically be used to generate evidence that meets regulatory requirements. Evidence from these trials is specifically relevant when treatment options already exist for the disease under study and when the real-life situation, including extraneous factors, is expected to influence the treatment effect. The process of randomisation in randomised controlled trials (RCTs) removes confounding by known and unknown factors. However, when a randomised control arm is not possible, an **external controlled arm (ECA)** can estimate the comparative treatment effect.

Typically, the external control arm uses data from past traditional clinical trials, but in some cases, RWD has been used as the basis for external controls. Using external controls has limitations, including difficulties in reliably selecting a comparable population because of potential changes in medical practice, lack of standardised diagnostic criteria or equivalent outcome measures, and variability in follow-up procedures. Collecting RWD on patients currently receiving other treatments, together with statistical methods such as propensity scoring, could improve the quality of the external control data when randomisation may not be feasible or ethical, provided there is adequate detail to capture relevant covariates.¹⁷

Hybrid prospective designs (e.g. concurrently randomised control as well as external control) allow traditional randomised controlled trials to be integrated with pragmatic design aspects to collect real-world data on patients. This design preserves the benefit of randomisation and provides real-world outcome data while potentially accelerating product development and lowering the cost of data collection and patient follow-up.

2. Planning and Study Set-Up

The objective of conducting clinical studies is to distinguish the effect of the drug from other influences, such as spontaneous change in the course of the disease, placebo effect, or biased observation. When relying on a non-interventional study (e.g. EHR data generated during routine clinical care analysed using a cohort study design), the inference(s) drawn may be incorrect if based on estimates that are affected by (1) confounding (e.g. due to noncomparable treatment groups) or (2) other forms of bias. Accordingly, before choosing a non-interventional study design for a study intended to support regulatory decisions regarding product safety and effectiveness, sponsors and

researchers should consider how likely it is that such a study design and its conduct will be able to distinguish a true treatment effect from other influences. Sponsors should address commonly encountered challenges when considering using a non-interventional study for regulatory decision-making.^{2,1}

Although statistical programmers are generally not involved in the early stages of clinical studies, it's still recommended to be aware of the processes upstream, especially if a regulatory submission is planned.^{2,1}

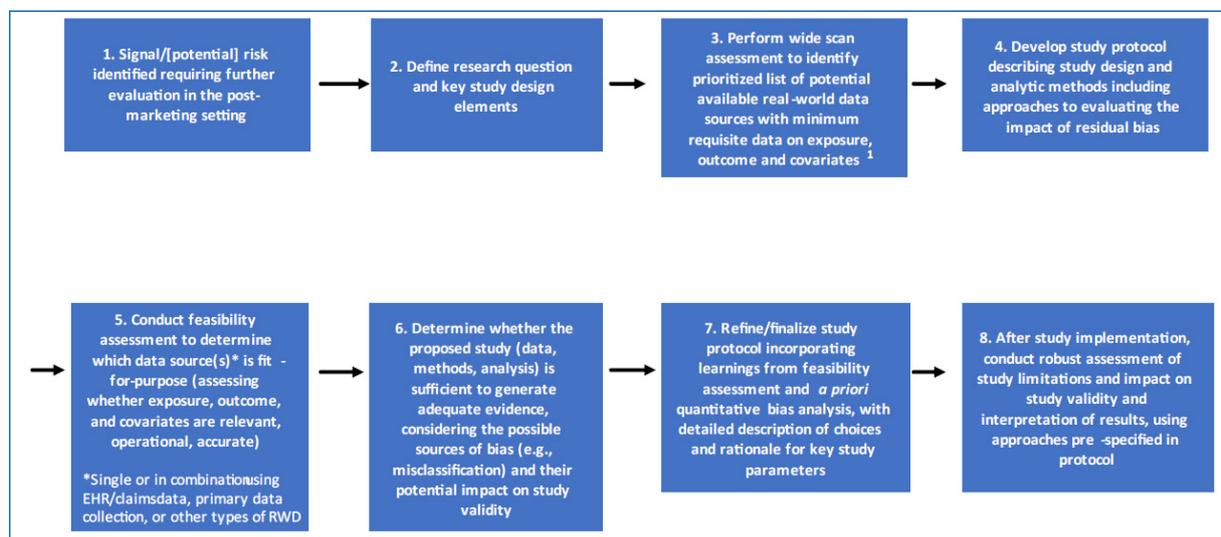


Figure 2.1: A framework for developing adequate evidence using fit-for-purpose real-world data to address regulatory questions on drug safety. Gatto NM et al. *The Structured Process to Identify Fit-for-Purpose Data: A Data Feasibility Assessment Framework*

2.1. Framing of the Research Question

The research question is a concise statement of the study purpose and the prespecified hypotheses to be tested. The purpose of the study may also be to generate hypotheses for future research. It is generally assumed that the research question has been defined and agreed upon ahead of statistical programming engagement. In the protocol, researchers should document decisions about the study design and the types of data required/available. Careful formulation of the research question will highlight unknowns that will need to be addressed through information derived from the feasibility assessment, and this information may further refine the question and drive protocol development.^{2,2}

2.2. Feasibility Assessment

A feasibility assessment is a systematic process to identify fit-for-purpose data to address a specific research question. When conducting a study-level feasibility assessment, a key goal is to describe and compare the reliability and relevance of the data sources assessed for the research question. Statistical programmers may be involved in this assessment. Feasibility assessments should be structured in at least two phases:

- An initial scan to determine whether the available data sources will suffice and to narrow down data source options
- A subsequent, more comprehensive feasibility assessment of the narrowed-down data sources.

In the early stages of designing a non-interventional study, sponsors should discuss with the regulators the expectations regarding access to patient-level or analytic datasets. Sponsors should obtain any agreements relevant to patient-level/analytic data required for submission. Submission of the feasibility assessment report can either be a standalone document, an

annex to the protocol, or used as context for design decisions in the protocol. The final approach should comply with applicable regulatory requirements. Detailed frameworks, templates and checklists for conducting feasibility assessments are available in scientific publications.^{2,3}

2.3. Engagement with Regulatory Health Authorities

Considering the evolving and diverse regulatory frameworks, early engagement with regulatory agencies is highly recommended. Before conducting non-interventional studies, sponsors are advised (but this is fast becoming a requirement) to submit the draft protocol and Statistical Analysis Plan (SAP) to the relevant RHA. For the FDA, sponsors should finalise the study protocol, including the research question of interest and rationale for the study design, before initiating study conduct.

The FDA strongly encourages sponsors to engage with the agency in the early stages of designing a non-interventional study and to provide sufficient information for clarifying expectations relating to the design and proposed conduct of their study. Although detailed information on every attribute described may not be available or feasible at the time of early engagement with the FDA, successful proposals for non-interventional study designs should satisfactorily address key attributes such as summary of the proposed approach, study design, data sources, and analytic approach, as applicable.

When the available data sources do not support proposals that satisfactorily address these attributes, alternative study designs should be considered.

2.4. Protocol Development and Review

Development of the protocol (and the SAP) occurs with alignment with the research question. The feasibility analysis will

guide the development of the protocol and facilitate discussion with regulatory health authorities, HTA bodies, and other parties. A major factor in bolstering confidence in RWE studies and ultimately producing RWE strong enough for regulatory decision-making is selecting fit-for-purpose data before finalising the protocol.

The protocol should contain details of the study design, data sources and data types, target population, exposure, outcomes, covariates, and the proposed analytical approaches. The study protocol and the SAP should specify the data provenance (curation and transformation procedures used throughout the data life cycle) and describe how these procedures could affect data integrity and the overall validity of the study.

Although statistical programmers are not part of the Protocol Review Committee (PRC), a statistical programming representative needs to be involved in the review of protocols describing RWE studies that are planned to be submitted. Statistical programmers should focus their review on areas that could impact the ability to interpret, transform, analyse or pool data.

2.5. Real-World Data Checklist

At the start-up, statistical programmers who are involved in the study can produce a checklist for internal use. Although the primary user and owner of this checklist is anticipated to be the statistical programming group at sponsor and CRO organisations, groups such as biostatistics, data management, project management and the RWE group may benefit.

The validity of the content of this checklist can be revisited after the clinical study reaches milestones such as finalisation of the protocol/the SAP and other study-specific milestones.^{2,4}

2.5.1. Purpose of the Real-World Data Checklist

Since the use of RWD is increasing in both interventional and observational clinical studies, users are recommended to apply this checklist in studies that include elements of RWD/RWE. Any such use of RWD in a clinical study can have implications in terms of how the data can be submitted to regulatory bodies. Please refer to Appendix 1.1 and example cases in Appendix 1.2 for how the checklist can be used. The checklist serves the following purposes:

- 1) It captures initial information about the nature of the real-world evidence study, with insights into statistical analysis and reporting functions. It helps the statistical programming group get a thorough view of the study in reference to the use of real-world data.
- 2) It helps the user understand real-world elements of the data and submission-related aspects – such as sources of data and effort involved – in meeting regulatory submission requirements.
- 3) It serves as preliminary guidance for initiating dialogue with stakeholders within the organisation on study planning, statistical analysis, and reporting activities.
- 4) The outcome of the checklist can help the user with decision-making for RWD submission planning. Examples can be found in Appendix 1.2.

3. RWD Study Council

As clinical research continues to add RWD to study designs, multistakeholder teams must evolve as well. The existing RCT clinical team is no longer sufficient to properly navigate RWD considerations. The proposed name for this new team is the RWD Study Council. Programmers will find themselves doing various titles depending on the organisation structure, but there are general principles that can be highlighted as RWD Study Councils go from sparse to common in the coming years.

3.1. Purpose

This chapter will unpack how an RWD Study Council functions, including general team member roles and how a responsibility assignment matrix can assist in visibility across the team, with an emphasis on details relevant to programmers. Therefore, given the scope of this white paper for programmers, other roles will need to expand beyond what is presented here.

3.2. RACI Matrix

One of the most common responsibility assignment matrices is the RACI matrix, which will be used here as a recommended practice for task clarity. Numerous sources describe the RACI acronym. We have selected one from a Forbes Advisor article:

“**Responsible** designates the task as assigned directly to this person (or group of people). The person responsible is the one who does the work to complete the task or create the delivery. Every task should have at least one responsible person and could have several.”

“The **accountable** person in the RACI equation delegates and reviews the work involved in a project. Their job is to make sure the person responsible, or the team, knows the expectations of the project and completes work on time. Every task should have only one accountable person and no more.”

“**Consulted** people to provide input and feedback on the work being done in a project. They have a stake in the outcomes of a project because it could affect their current or future work.”

“**Informed** folks need to be looped into the progress of a project but not consulted or overwhelmed with the details of every task. They need to know what’s going on because it could affect their work, but they’re not decision-makers in the process.”

The RACI matrix is a proven resource and will be especially helpful as RWD Study Councils get up and running. A general RACI matrix template is proposed below, including common contributors/roles and activities, followed by company-specific examples.

3.3 Contributors/Roles

The chart in Figure 3.1 shows possible role dependencies. The roles in the top level are included in Table 3.1.

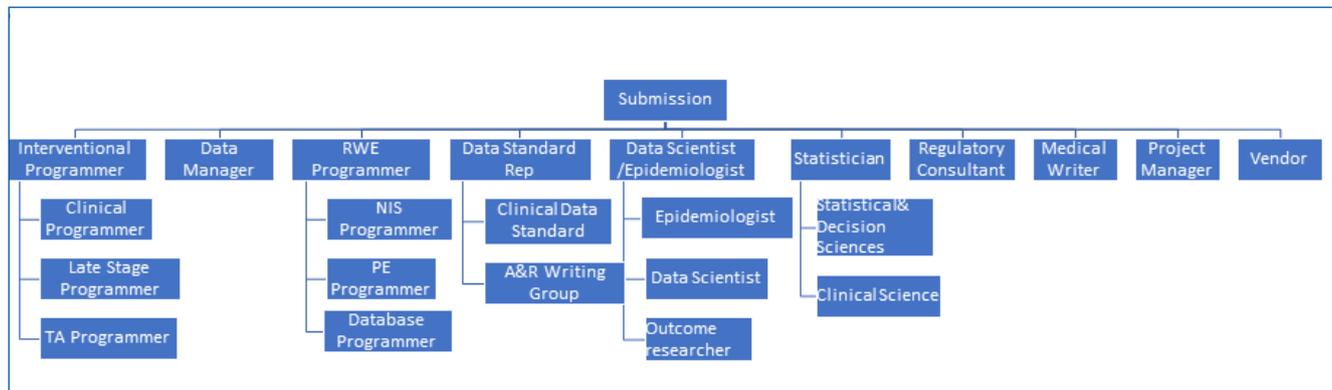


Figure 3.1. RWD – Hierarchy of Study Roles

TA – Therapeutic Area, PE – Pharmacoepidemiology, NIS – Non-Interventional Study, A&R – Analysis & Reporting

The roles for participants in RWD are included as separate columns in Table 3.1. Each company may use different naming conventions based on the department structure and related internal processes. It is important to note that both RCT and RWE programmers are required to support RWE studies. The individual rows present the main tasks that can be considered for the submission preparation process.

There are four main study phases listed: Study Initiation (kick-off meetings, project management tasks, etc.), Study Execution (data management and data handling rules/standards, etc.), Submission Preparation (regulatory agency requirements etc.), and Close-Out (including final submission packages and sharing data).

3.4. General RACI Table

Study Phase	Activity/Roles	Interventional Programmer	Data Manager	RWE Programmer	Data Standard Rep	Data Scientist /Epidemiologist	Statistician	Regulatory Consultant	Medical Writer	Project Manager	Vendor Partner
Study Initiation	Planning										
	Complete Roles and Responsibilities (RACI)										
	Project Management/ Kick-off Meeting										
	Identifying the RWD sources										
Study Execution	Feasibility Assessments										
	Data Handling										
	Data Transfer/Data Share Agreements										
	Considerations										
	Data Standardization Plan										
	Coding Versions										
	Milestone dates agreement										
	Data Analytics										
	SDTMs/ADaMs and Metadata										
	Tables/Listings/Figures										
Submission Preparation	SDTMs/ADaMs										
	Submission Package Review										
	Submission Requirements										
Close-out	Regulatory Reviews										
	Final Submission Package										
	Sign off final documents										
Submit to regulatory agency											

Table 3.1. RACI – RWD – Study Roles and Activities

Table 3.1 serves as an RACI matrix template. It provides a high level of categories for individual contributors and activities involved in RWE submission. Below are two approaches to adapting an RACI template based on the company’s requirements and/or organisational structure.

3.4.1. RACI Table Examples

Company A

There is no separation between RWD and RCT programmers in company A. Clinical Programming (CP) and Statistical Programming (SP) constitute the programming group that deals with both types. For full RWE studies, the Data Sciences (DS) group is the driver for analysis (equivalent to stats in RCT), and they come to CP and SP when submission is planned.

Activity	CDS&T	DM	RM-CM	CP	SP	RMW	SDS	DS/Epi	CS	GTL/ GPL
Identifying the RWD sources	I	I		I	I		R, A	R, A	I	I
Feasibility Assessment: mini analysis plan	I	I	I	I	C	C	C, R	R, A	I	
Feasibility Assessment: Determining Data Relevancy		I	I	I	C	C	C, R	R, A	I	I
Feasibility Assessment: Data Reliability analysis		R	I	R	R		R, A	R, A	I	I
Feasibility Assessment and Report		C	C	C	C, R	R	R, A	R, A	I	I
Protocol (using appropriate template)	C	C	C	C	C	R, A	C	C	C	I
Vendor Engagement	C	C, R	I	C	C	I	C	C, R	I	R
Data Ingestion Source Data Flows		R, A		C	C		I	I	I	
Coding	C	R, A		C	C		I	I	I	
Database Releases		R, A		C	C		C	C	C	
Creation of Mapping File		I		R, A	C		C	C	I	
Creation of SDTMs or CDM	C	R, A		R, A	C		I	I	I	
Creation of SAP		I		C	C		R, A	I	I	
Creation of DPS	C	I		I	R, A	C	R	C, I	C	
Creation of ADaMs	C	I		I	R, A		C	I	I	
Data linkage from more than one Databases	C	R		C	R, A	I	R	R	R	
Data Integrations		I		I	R, A		C	C	C	
Standards used	C	I		R, A	R, A		R, A	R, A	I	
Determination of waiver requirement for SDTMs	C	I		R, A	R		C	I	I	
Determination of waiver requirement for ADaMs		I		I	R, A		C	I	I	
Study Data Standardization Plan	R	C		C	C					
Privacy / Legal Considerations	C	I		I	I		R, I	R, I	R, I	R
Data De-identification	R									

Note: Abbreviations: CDS&T = Clinical Data Standards & Transparency, DM = Data Management, CP = Clinical Programming, SP = Statistical Programming, RMW = Regulatory Medical Writing, SDS = Statistics and Decision Sciences, DS = Data Sciences, Epi = Epidemiology, CS = Clinical Sciences, GTL = Global Trial Lead, GPL = Global Program Lead.

Company B

There is a separation between RWD and RCT programmers in company B. RWD programmers fall under two roles: Therapeutic Area (TA) and Regulatory RWE Programming (RRP). TA programmers are given an SAP to programme against. They do the actual regulatory study. RRP is there to guide the TA programmers through the data standards and train on the templates, Excel spec, P21, reviewer's guides and the SAS macros to make sure they are fulfilling regulatory and company B requirements and processes.

Activity	TA Programming Lead	TA Programmer(s)	TA Epidemiology Lead	RRP Lead	RRP Programmer(s)	PMO	Cross Functional Partners: RCT Programmer(s)	Consultants: SDTM / ADaM Standards
Study Initiation								
Submit RRP Parent DAC Request for Regulatory Study	I		R	I		I		
Submit PMO Request	I		I	R		I		
Schedule PMO Kick-off Meeting	I	I	I	I	I	R		
Study Setup								
Create MS Teams for Regulatory Study	C	I	I	C	I	R	I	
Create Smartsheet Risk Register	I		I	I		R		
Create Project Status Dashboard Slide	I		I	I		R		
Create Question and Decision Log	I		I	I		R		
Assess Data Submission Needs (e.g. SDTM, ADaM, P21E license, etc)	R	I	C	R	I		C	C
Complete Roles and Responsibilities Form	R	I	R	C	I			
Complete Data and Folder Access Form	R	I	I	C	I			
Submit RWD Engineering Child DAC Request	R			I		I		
Raw Data Access and Preparation	I	I	I	I	I	I		
Create Regulatory Project Folder	I	I	I	I	I			
Populate Regulatory Project Folder (templates, etc)	I	I	I	A	R			
Assess Clean Room Needs	R		R	R				
Open and Setup Clean Room	I	I	I	I	I			
Submit Approved Clean Room Requests	C	I	R	C	I			
Fulfill Clean Room Requests	A	R	C	C	C			
Study Request Review								
Study Analysis								
Ad Hoc Analysis								
Study Sign-off								

Note: Abbreviations: TA = Therapeutic Area, RRP = Regulatory RWE Programming, PMO = Project Management, RCT = Randomised Controlled Trial, DAC = Data Analytics Centre.

3.4.2. Conclusion

As the need for RWD Study Councils continues to grow, team members need clarity on their roles and responsibilities. The RACI matrix provides this clarity and can be adapted to organisational structure. RWD programmers can offer this solution to their study teams if it does not exist.

4. Data Provenance, Data Ethics and Data Privacy

4.1 Data Provenance

Data provenance refers to the complete documentation of a dataset's lineage – including its origin, collection context, transformation processes and any subsequent modifications. As a form of metadata, provenance serves as an audit trail that enables verification of a dataset's authenticity, supports reproducibility, and ensures transparency and accountability in scientific research⁴¹ Provenance plays a critical role in validating the integrity of clinical and real-world data and is essential for evidence generation in regulatory and health technology assessments. In parallel, data privacy refers to the protection of sensitive personal information from unauthorised access, use or disclosure. In today's data-driven landscape, particularly in the context of clinical trials, digital health and genomics, data privacy is challenged by the need for openness, traceability and interoperability.

While data provenance enhances transparency and trust, it can also pose privacy risks, especially when provenance metadata includes sensitive contextual details that may

inadvertently reveal personal or institutional identities.⁴² Thus, there is a critical need to balance transparency with confidentiality. Provenance frameworks must be designed to respect data minimisation principles and incorporate privacy-preserving techniques, such as de-identification, anonymisation, tokenisation, hashing and encryption, to ensure sensitive information is not exposed.⁴³ This balancing act becomes particularly important in clinical research, where enabling data sharing and leveraging emerging technologies such as artificial intelligence, distributed learning and cloud-based analytics introduces new risks – including data theft, re-identification and unauthorised manipulation. De-identification remains the most common privacy-enhancing method, yet its effectiveness varies across regulatory contexts. For instance, GDPR requires anonymised data to be irreversibly de-linked from individuals, while other frameworks, such as the US HIPAA Safe Harbor standard, apply more pragmatic thresholds based on re-identification risk.^{44, 45}

4.2 Data Ethics

Data ethics broadly refers to the moral principles guiding how individuals, organisations and institutions collect, manage, protect and use data. It encompasses core values such as fairness, transparency, accountability and respect for privacy, emphasising that data governance should serve both individual rights and the public good.⁴⁶ Ethical data practices are particularly critical in healthcare, where the misuse or overuse of sensitive health information can lead to discrimination, stigmatisation, and erosion of public trust.

4.3. Data Privacy

4.3.1 Adhering to Data Privacy Regulations

Complying with data privacy regulations is essential for safeguarding individual data, fostering stakeholder trust, and mitigating legal and financial risks. Key legislative frameworks – such as the General Data Protection Regulation (GDPR), the California Consumer Privacy Act (CCPA) and the Health Insurance Portability and Accountability Act (HIPAA) – establish stringent requirements for data collection, processing and disclosure. Adhering to these regulations necessitates a comprehensive approach that includes respecting data subject rights (e.g. access, rectification and erasure), conducting Data Protection Impact Assessments (DPIAs), implementing timely breach notification protocols, and ensuring accountability across third-party data processing arrangements.^{4,7,4,8,4,9,4,10,4,11}

- **General Data Protection Regulation (GDPR) (2018)** – A comprehensive data protection law enacted by the European Union, GDPR governs the collection, processing and storage of personal data. It strengthens individual privacy rights, mandates lawful data handling practices, and applies extraterritorially to any entity processing EU citizens' data.^{4,12}
- **Personal Information Protection and Electronic Documents Act (PIPEDA) (2000)** – Canada's federal privacy law for private-sector organisations, PIPEDA governs how businesses collect, use and disclose personal information during commercial activities. It emphasises individual consent, data access rights and the requirement to safeguard personal information.^{4,13}
- **The Privacy Rule of the Health Insurance Portability and Accountability Act (HIPAA)** – A foundational US regulation that protects individuals' medical records and other personal health information (PHI). The HIPAA Privacy Rule sets national standards for the use and disclosure of PHI by covered entities and gives patients rights to access, amend and request restrictions on the use of their health data.^{4,14}

Growing global emphasis on transparency has led regulatory agencies to implement policies that govern the public disclosure of clinical trial data. The following key initiatives illustrate how major regulators have shaped the data sharing landscape:

- **EMA Policy 0043 (2010) and Policy 0070 (2015)**
The European Medicines Agency (EMA) introduced Policy 0043 to enable public access to documents held by the agency, including regulatory submissions. Policy 0070, implemented in 2015, advanced transparency by proactively publishing clinical study reports submitted in marketing authorisation applications, while protecting personal and commercially confidential information.^{4,15}
- **Health Canada Public Release of Clinical Information (2019)**
Implemented in 2019, Health Canada's PRCI initiative mandates the proactive publication of anonymised clinical information from drug and medical device submissions following final regulatory decisions. This policy aims to enhance transparency, support independent research and align with international best practices, while ensuring the protection of personal and confidential business information.^{4,16}

- **EU Clinical Trial Regulation (2022)**

Effective from 31 January 2022, the EU Clinical Trial Regulation harmonises the assessment and supervision processes for clinical trials across EU Member States. It introduces the Clinical Trials Information System (CTIS), a centralised portal facilitating single-submission applications for multinational trials, thereby enhancing efficiency, transparency and participant safety.^{4,17}

4.3.2. Statistical Disclosure Methods

Statistical disclosure control (SDC) encompasses a set of techniques designed to minimise the risk of re-identifying individuals from statistical outputs while preserving the analytical value of the released data. These methods are particularly critical in the dissemination of aggregate statistics, microdata, and synthetic datasets derived from sensitive sources such as clinical trials, electronic health records or administrative databases.^{4,18}

At its core, SDC seeks to uphold data confidentiality and ensure information disseminated to researchers, policymakers or the public complies with ethical principles and legal frameworks such as GDPR, HIPAA and national statistical legislation. SDC techniques are used not only to prevent direct identification but also to mitigate the risks of inferential disclosure, where individuals can be indirectly identified through combinations of attributes or auxiliary information.^{4,19}

Types of Statistical Disclosure Control Techniques

SDC methods can be broadly categorised into two types:

Input SDC (applied to microdata): These methods are used before statistical analysis and include:

- Data swapping
- Top- and bottom-coding
- Noise addition
- Data suppression
- Record aggregation or micro aggregation.

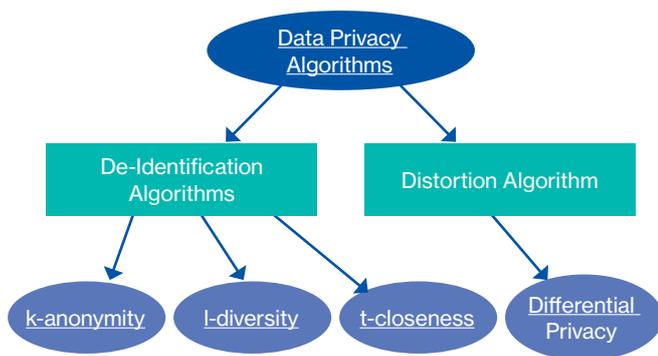
Output SDC (applied to analytical results): These techniques are applied to summary statistics and outputs such as tables or regression coefficients and include:

- Threshold rules (e.g. minimum cell counts)
- Rounding and perturbation of cell values
- Suppression of sensitive cells in tabular outputs
- Dominance rules to prevent disproportionate influence of small subgroups.

The challenge with implementing SDC lies in the trade-off between disclosure risk and data use. Excessive masking can render datasets analytically useless, while insufficient masking can compromise individual privacy (Garfinkel, 2015). Therefore, selecting the appropriate SDC method requires a contextual risk assessment, often guided by formal risk-utility frameworks or privacy risk metrics.^{4,20}

4.3.3. Data Privacy Methods

Data privacy techniques employed in clinical research and real-world data analysis can broadly be categorised into two major methodological classes: distortion-based methods and de-identification-based methods. Distortion methods, such as differential privacy, data perturbation, and noise addition, deliberately modify data values to obscure individual-level information while preserving statistical utility.⁴¹⁶ These approaches offer strong formal privacy guarantees, especially in large-scale data environments, but often require careful tuning to balance privacy with analytical accuracy.⁴¹⁴ In contrast, de-identification methods focus on removing or masking direct and indirect identifiers (e.g. names, birth dates, ZIP codes, and rare diagnoses) to reduce re-identification risk. Techniques such as pseudonymisation, k-anonymity, l-diversity and data generalisation fall under this category.⁴²¹ While de-identification is widely used in regulatory compliance (e.g. under HIPAA and GDPR), it can be vulnerable to linkage attacks when auxiliary information is available.⁴²² In practice, a hybrid approach – combining statistical distortion with structural de-identification – is often necessary to meet both privacy protection and data utility requirements.⁴³



Source: Biswas, S., Fole, A., Khare, N. et al. (2023). *Enhancing correlated big data privacy using differential privacy and machine learning*. *J Big Data* 10, 30.

- *The k-anonymity* method is based on obscuring individual identities in a dataset by grouping at least k similar individuals and suppressing identifying attributes. This reduces the risk of disclosing personal information about any individual within the group.⁴²¹
- *The l-diversity* method is an extension of the k-anonymity method, which enhances privacy in datasets by reducing the granularity of data representation and providing protection against attribute disclosure. For example, each group of k individuals should contain at least l different medical conditions to prevent inference about sensitive attributes.⁴¹⁷
- *The t-closeness* model is an extension of the l-diversity framework, which further improves privacy protection by considering the distribution of sensitive attribute values within each anonymised group. It ensures the distance between the distribution of a sensitive attribute in any group and the distribution of that attribute in the overall dataset does not exceed a threshold t , thereby limiting the risk of attribute disclosure.⁴²³

- *The Differential Privacy (DP)* is a formal privacy framework introduced in theoretical computer science that ensures the risk of identifying an individual remains low, even with access to auxiliary information. An algorithm satisfies DP if the presence or absence of a single individual's data does not significantly alter the output of a query. It provides quantifiable privacy guarantees and supports data sharing through six key properties, including resistance to linkage attacks, composition, and post-processing immunity.^{415,416}

The integrity and trustworthiness of clinical and real-world evidence hinge on the seamless integration of data provenance, ethics and data privacy. Data provenance ensures transparency and traceability of data from its origin through all stages of processing and analysis. This traceability is essential for upholding ethical standards, including informed consent, fairness, and respect for people. In turn, data privacy safeguards individual rights by applying rigorous protections to sensitive information, reinforcing both ethical obligations and the credibility of the data's lineage. Together, these pillars form an interconnected framework that supports responsible data stewardship, enabling high-quality, ethically sound, and privacy-compliant research.⁴²⁴

5. Vendor Engagement

Vendor engagement (VE) in the context of RWD and EHR (electronic medical records) is the collaboration between healthcare organisations (such as hospitals, clinics and medical practices) and vendors that provide EHR systems. Vendor engagement is essential for implementing EHRs and their ongoing use in healthcare organisations. This partnership ensures EHR systems meet the needs of healthcare providers and improve patient care.

Here are some EHR vendors:

- Epic (<https://www.epic.com>)
- Oracle Health EHR (<https://www.oracle.com/health/clinical-suite/electronic-health-record/>)
- athenahealth (<https://www.athenahealth.com>)
- eClinicalWorks (<https://www.eclinicalworks.com>)
- Veradigm (<https://veradigm.com/>)
- MEDITECH (<https://ehr.meditech.com/>)

5.1. Selection

Healthcare organisations engage with EHR vendors to select EHR systems that best fit their needs. This involves evaluating the products, features, pricing and support services.^{5.1,5.2,5.3,5.4,5.5}

5.2. Implementation

This involves data migration, software installation, customisation to meet workflow requirements, and staff training. The customisation could be to support specialised clinical workflows, integrate with other software solutions, or implement additional features and modules.

5.3. After Implementation

Vendor engagement continues after the EHR system has been implemented. Healthcare organisations rely on vendors for ongoing technical support, software updates and maintenance

to ensure the systems function smoothly and comply with regulatory requirements.

Feedback and collaboration between healthcare organisations and the vendor on usability, functionality and performance of the EHR system can inform updates and improvements.

5.4. Regulatory Compliance

Regulatory compliance in the context of EHRs typically includes adhering to laws and regulations aimed at protecting patient privacy, ensuring data security and promoting the interoperability of health information. Key regulations healthcare organisations and EHR vendors need to comply with include:

- The Health Insurance Portability and Accountability Act (HIPAA)
- The HITECH Act
- The Promoting Interoperability Program
- Regulations by the FDA and others.

5.5. Case Studies

5.5.1. Case Study 1: Vendor Engagement in a Real-World Data Study

Background

In a real-world evidence trial focused on analysing dual combination therapies used in the treatment of hypertension within a multinational cohort, vendor engagement played a crucial role. This case study illustrates key aspects of vendor engagement in RWD and RWE studies without promoting any specific organisation.

Importance of Vendor Engagement

Vendor engagement is important for several reasons:

- Access to diverse and large-scale data: Vendors provide access to extensive healthcare data across regions and practice settings, which is essential for conducting large-scale, multinational research studies.
- Standardisation of data: Collaborating with vendors enables the use of standardised data models, ensuring consistency and comparability across studies and databases.
- Regulatory compliance: Vendors often have established protocols for managing data in compliance with regulatory requirements, ensuring patient privacy and data security.
- Technical support and expertise: Vendors offer technical support and expertise in data management and analysis, facilitating more efficient and effective research processes.
- Application of advanced technologies: Vendors may apply AI solutions across the product life cycle, bringing precision, speed and scale to stages of the research process.

Implementation

Vendor engagement typically begins at the early stages of a research project and continues throughout the study to ensure ongoing access to data, technical support and compliance with emerging regulatory requirements.

- Initiation and planning: The research team identifies the need for external data or analytics capabilities and reaches out to potential vendors.
- Negotiation and agreement: Terms of access, use of data, costs, and compliance with privacy laws are negotiated and

formalised through data use agreements or contracts.

- Data access and standardisation: The vendor provides access to data, often involving data transformation to fit standard models.
- Ongoing support and collaboration: The vendor offers technical support, training and consultation throughout the project.
- Compliance and ethical oversight: Both parties ensure the project complies with ethical standards and regulatory requirements through established frameworks.

Monitoring and Assessment

Effective monitoring and assessment of vendor engagement ensures the collaboration meets research objectives and adheres to agreed standards and regulations. This process involves:

- Performance metrics: Establishing clear performance metrics related to data quality, timeliness of data delivery, and adhering to the project timeline and budget.
- Regular updates and meetings: Scheduling regular meetings with the vendor to review project progress, discuss challenges and adjust plans as necessary.
- Compliance checks: Regularly reviewing processes and data handling practices to comply with regulatory requirements and data security standards.
- Technical support and issue resolution: Monitoring how responsive and effective the vendor is in providing technical support and resolving issues.
- Data quality and utility: Assessing the quality of the data provided by the vendor and how useful it is for meeting research objectives.
- Satisfaction surveys: Conducting surveys or interviews with the research team to gather feedback on their experience of working with the vendor.
- Cost-effectiveness: Analysing the cost-effectiveness of the engagement, considering the value of the services provided in relation to the cost and impact on the project's budget.

5.5.2. Case Study 2: Vendor Engagement Process for the FDA-Approved RWE Study

This case study looks at how vendor engagement is practised in a typical pharmaceutical organisation, specifically vendor selection strategies, RWD collection and management, data standardisation across sources, and the pivotal role of programming and validation in ensuring data quality and analysis readiness.

Vendor Selection and Data Collection

Vendor selection for data sourcing depends heavily on the nature of the study analysis. A multidisciplinary team, including members from medical affairs, real-world data science, and procurement, assesses and selects vendors. This process is driven by study objectives, data availability from licensed sources, and the need for external data acquisition versus in-house capabilities.

In studies where patient data is collected for research purposes, particularly those with less stringent budget constraints, the company has encountered challenges in consistently capturing and abstracting the required data. This variability is due to the absence of a universal data model such as CDISC standards. However, efforts are underway to improve data capture frequency and standardisation. Currently, each study may rely on unique data models, necessitating tailored programs for

standardising and analysing the collected data.

Data Access and Standardisation

In the real-world data landscape, standardisation levels differ significantly from those in clinical trials. Clinical trials adhere to stringent standards primarily for regulatory acceptance and data integrity assurance. In contrast, claims databases prioritise data standardisation to facilitate billing processes for healthcare providers. Electronic medical records (EMRs) have their own set of standardisation principles, dictating what data should be collected, though specific table structures and field formats vary widely.

The data dictionary serves as a critical tool for understanding and structuring EMR data. It typically includes structured data elements such as table names, field descriptions, data types, lengths, and formats. This dictionary aids data interpretation and ensures compliance with privacy regulations. Meanwhile, data models such as OMOP and Sentinel's common data model play crucial roles in standardising and harmonising data from multiple sources, streamlining the analysis process.

Developing programming specifications for data analysis is a nuanced process, tailored to the data source's characteristics. Protocols and Statistical Analysis Plans (SAPs) guide this process, outlining analytical objectives and methodologies. Flexibility is essential in real-world data analysis, as unforeseen insights often necessitate adaptive approaches to data interpretation and programming.

The operational manual or user guide for databases provides comprehensive insights into database use, including enrolment, expenditures, demographics and clinical procedures. Each database within the system specifies the tables and content available. For example, Medicaid databases encompass diverse information about patient care coverage, services and clinical outcomes, organised into data tables based on predefined client specifications.

Monitoring and Validation

External data studies require meticulous oversight by cross-functional stakeholders to monitor progress, identify anomalies, and address data-related issues promptly. Routine meetings ensure alignment with study milestones and facilitate timely interventions to maintain data integrity.

Data validation tools such as SAS, R and Python play a pivotal role in ensuring data accuracy and quality. They enable early detection of errors and adherence to predefined validation rules and edit checks. These tools are instrumental in preparing clean datasets for analysis, verifying data completeness, and identifying any anomalies or missing values.

Summary

In summary, using real-world data (RWD) in clinical research and healthcare analytics involves navigating a landscape of diverse data sources and varying standardisation requirements, and vendor engagement is an important part of having high-quality healthcare data. The challenges of data collection, particularly in studies with budget constraints, highlight the importance of developing tailored data models and programming specifications to standardise and analyse data effectively. Vendor selection is a strategic process driven by study objectives and the availability

of licensed data sources. Data dictionaries and standardisations are crucial in structuring and interpreting electronic medical record (EMR) data for compliance with privacy regulations. Beyond data collection and standardisation, effective monitoring and assessment of vendor engagements, exemplified by collaborations with CROs in research projects, is critical. This involves establishing clear performance metrics relating to data quality and timeliness, scheduling regular updates and compliance checks, and evaluating data utility, timeliness of delivery and cost-effectiveness. Such practices ensure research collaborations meet objectives, adhere to standards and regulations, and ultimately enhance the value and reliability of real-world data for informed decision-making in healthcare.

6. Fit for Purpose Assessment

Before using real-world data (RWD) in research, sponsors must assess whether the data is appropriate for the research question. This 'fit for purpose' assessment or 'feasibility analysis' ensures RWD is both relevant and reliable for use in a clinical investigation. Relevance refers to the availability of data for key study variables – such as exposure, covariates and outcomes – and whether it includes enough representative subjects. Reliability refers to the accuracy, completeness and traceability of the data. Without this, conclusions drawn from the study may be flawed.

Fit for purpose assessment is a critical step for determining whether RWD is suitable for study. Researchers need to evaluate the data sources for potential biases, gaps and limitations. For example, while injectable or intravenous medications are often captured accurately in RWD, oral or inhaled drugs may be harder to track. This means a data source that works well for one study might not be appropriate for another, even if the studies and/or data sources appear to be similar on the surface.

After identifying the appropriate study design, sponsors evaluate RWD sources to determine whether they are reliable and relevant enough to answer the research question. However, unlike randomised controlled trials (RCTs), which minimise bias through randomisation, studies based on RWD require alternative strategies to control for observational bias when evaluating sources of RWD to generate valid real-world evidence (RWE) within the context of a clinical investigation.

6.1. Scenario: Designing a Study Using RWD

Suppose researchers are evaluating the efficacy of a new medication based on RWD sources. Before collecting data, they define a hypothetical target trial (HTT) based on a double-blind, randomised controlled trial. The intent to emulate the HTT forces sponsors to account for known confounding factors such as patient demographics, medical history and treatment protocols, which would otherwise be independently and identically distributed across treatment arms due to randomisation. For example, if treatment arms have imbalanced age groups, the results could skew the perceived effectiveness of the intervention.

6.2. Defining a Hypothetical Target Trial (HTT)

Defining an HTT helps researchers anticipate sources of bias that could affect the study results and helps programmers identify those data sources. Both clinical and statistical

expertise are needed to clearly articulate the key design elements for evaluating RWD sources to determine whether they are fit for purpose:

- 1. Estimator**
 - o Define the statistical methods used to test the study's null hypothesis.
- 2. Treatment Group 'Assignment'**
 - o Specify who will receive the treatment and who will serve as the control group.
- 3. Time Zero**
 - o Determine the reference point from which changes or outcomes will be compared.
- 4. Length and Frequency of Follow-Up**
 - o Set the duration of participant observation and how often data will be collected during that period.
- 5. Sample Size**
 - o Ensure the study has enough participants to generate reliable results.
- 6. Inclusion/Exclusion Criteria**
 - o Clearly define who is eligible to participate in the study and who is not, based on the research objectives.
- 7. Threats to Validity**
 - o Identify any factors that could affect the accuracy of the study, and devise strategies to address them.
- 8. Secondary Outcomes**
 - o List additional metrics beyond the main study goal, define how they will be measured, and explain their relevance.
- 9. Key Subgroups**
 - o Identify subgroups within the study population that are critical for understanding treatment effects.
- 10. Confounding Variables**
 - o Recognise potential factors that might distort the relationship between treatment and outcome, and explain how these will be controlled.
- 11. Rationale for Confounder Selection**
 - o Justify why specific confounders are selected for control, with evidence supporting their inclusion.

Once the HTT is defined, sponsors can assess whether the data sources available in the real world adequately capture the necessary variables and whether those data points are reliable. If the data falls short, for example if adherence rates are not captured accurately, or if there are discrepancies in how treatments are documented, sponsors may need to revise their study design or seek additional data.

This iterative process ensures the data is fit for the purpose of answering the research question of the clinical investigation, with all decisions to include or exclude RWD from analysis based on objective evidence. The documentation produced includes data definitions and study-specific data quality checks to make it easier for reviewers to replicate the study and assess its validity.

7. Analysis and Submission

7.1. Statistical Considerations

RWE studies and RCTs differ significantly in terms of study design, control over variables, and handling of biases and confounding. Both approaches have unique strengths and limitations, and researchers must carefully consider these factors when interpreting study findings and making evidence-based decisions in medicine.

Observational nature: RWE studies are typically observational in nature, meaning they involve analysing data collected from routine clinical practice, electronic health records (EHRs), administrative claims databases, or registries. Data is collected without researcher intervention or manipulation.

Lack of randomisation: Unlike RCTs, RWE studies do not involve random assignment of participants to different treatment groups. Patients receive treatments based on real-world clinical decisions made by healthcare providers.

Confounding and biases: In RWE studies, confounding is a major concern due to the lack of randomisation. Patients may differ systematically based on factors that influence both treatment assignment and outcomes. Controlling for all potential confounders is challenging, and residual confounding can bias results. In addition, RWE studies are susceptible to selection bias, information bias and confounding bias. Bias can arise due to differences in patient characteristics, incomplete or inaccurate data, and unmeasured variables affecting outcomes.

7.2. Confounding and Biases

Confounding occurs when the relationship between an independent variable (such as a treatment or exposure) and a dependent variable (an outcome) is distorted by an additional factor (a confounder) that is associated with both the independent variable and the outcome. In simpler terms, confounding arises when a third variable influences both the exposure and the outcome, making it difficult to accurately assess the true effect of the exposure on the outcome.

Selection bias: This bias occurs when the study sample is not representative of the target population, which leads to erroneous conclusions. For example, if a study on a new medication only includes participants who are younger and healthier, the findings may not apply to older or sicker individuals.

Information bias: This bias arises from errors in the measurement or assessment of exposure, outcome, or confounders. For instance, recall bias occurs when participants do not accurately remember past exposures or outcomes, which causes misleading associations.

Measurement bias: Similar to information bias, this occurs when there are errors in how variables are measured, particularly relevant in RWD, where data may not be collected with research purposes in mind and thus may lack standardisation.

Attrition bias: This can happen in longitudinal RWE studies if there is a loss of participants over time, and this loss is not random but related to the characteristics of the individuals or their treatment.

7.3. Methods to Address Confounding in RWD

Study design: Establish clear protocols for data collection and analysis before initiating the study. This helps minimise the risk of selective reporting and data-driven decisions. Whenever possible, include appropriate comparison groups to control for confounding factors. For example, use control groups or historical comparators to assess the impact of interventions or exposures.

Statistical adjustment: Use regression models (e.g. logistic regression, Cox proportional hazards models) that incorporate multiple covariates to adjust for potential confounding variables. Include relevant patient characteristics, disease severity and other factors known to influence outcomes.

Matching: Match treated and control subjects based on key covariates (e.g. age, sex, comorbidities) using simple matching or propensity score matching. This creates comparable groups and reduces the impact of confounding.

Stratification: Stratify the study population based on key confounders (e.g. age groups, disease severity) and conduct separate analyses within each stratum. This allows for examination of treatment effects within specific subgroups while controlling for confounding variables.

Sensitivity analysis: Conducting sensitivity analyses to assess how sensitive results are to changes in the method of confounder adjustment or the presence of unmeasured confounding can provide insight into the robustness of the study findings.

Post hoc analysis: Using post hoc techniques to adjust for variables that become apparent, as confounders only during the analysis can further refine the study outcomes.

Transparent reporting: Clearly describe the methods used for confounding adjustment, including details on covariate selection, model specification and handling of missing data. In addition, transparently report study limitations, including potential sources of bias and unmeasured confounders. Acknowledge the inherent limitations of observational data and discuss the implications for result interpretation.

In summary, by implementing these effective strategies in RWE studies, researchers can enhance the validity and reliability of findings, despite the inherent challenges of confounding and biases. However, it's important to recognise that no method can eliminate bias, and careful consideration of study design and data quality remains essential in generating meaningful real-world evidence.

7.4. Submission of RWD: Current Regulatory Landscape

The FDA will focus on three areas when evaluating RWD submitted in support of a marketing application: 1. Whether RWD is fit for use 2. Whether the study design can provide adequate evidence 3. Whether the study conducted meets regulatory requirements. These are described in the *Framework for FDA's Real World Evidence Program* published in 2018.⁷¹

The FDA published detailed guidance on data standards required when submitting RWD in support of a marketing application in *Data Standards for Drug and Biological Product*

Submissions Containing Real-World Data Guidance for Industry (December 2023).⁷² The key takeaway is that RWD and data from non-interventional study designs must be submitted using the standards documented in the FDA Data Standards Catalog. This means that RWD must be submitted using Clinical Data Interchange Standards Consortium (CDISC) standards; thus, sponsors need to convert real-world datasets into CDISC format when submitting this information in support of a marketing application. This guidance also provides advice on mapping RWD in CDISC format and considerations for the required documentation. While the agency acknowledges that its current catalog of standards does not necessarily reflect data derived from real-world sources, it has indicated that it is considering updates, including recommendations for mapping.

The landscape for submitting RWD is still evolving. FDA guidance documents indicate that the data standards and documentation required for submitting RWD are the same as for submitting RCT data. This presents a lot of challenges for the sponsors. Most of these challenges stem from the fact that RWD, by definition, is not collected under the supervision of a protocol and by research study staff. In addition, the business cases for using RWD and RCTs differ. RWD is typically not designed or collected for research purposes. As a result, real-world databases are organised and configured in a way that makes sense to healthcare providers, not researchers. Furthermore, these databases use data standards that contain different concepts and coding systems to those defined within CDISC and used in regulatory review.⁷³

As far as exchange standards are concerned, the required standard for data exchange for the FDA is SAS V5 transport file format. Although this has been the predominant format for submission of clinical data, it imposes several restrictions: limited data types, alphanumeric variable names, limited length for variable names, labels and widths, to name a few. With the breadth and depth of RWD data sources, the SAS V5 transport format is becoming a difficult choice for sponsors and thus not fit as a viable option for data exchange. While CDISC datasets are currently exchanged in SAS V5 transport format, CDISC is not inherently tied to this format. CDISC datasets can also be exchanged using XML, JSON or other file formats. Therefore, if the FDA and other regulatory authorities agree, CDISC datasets may be exchanged using another format.⁷⁴

7.5. How the Regulatory Submission Process for RCT Translates into RWD

Overview of the Clinical Study Data Reviewer's Guide (cSDRG)

Legacy tabulation datasets:

CDISC standard cSDRG can be created, which provides information about raw datasets and terminology that benefit from additional explanation beyond the Data Definitions document (Define-XML). Summary of SDTM conformance findings cannot always be included in this reviewer's guide, because the raw legacy data has not been converted to SDTM format. The standard cSDRG template can be used for filing. Details about XPORT files and example codes for importing the XPORT files into SAS or R need to be written in the Data Format and Import Information section.

Overview of the Analysis Data Reviewer's Guide (ADRG)

Legacy analysis datasets:

CDISC standard ADRG can be created, which provides information about analysis datasets and terminology that benefit from additional explanation beyond the Data Definitions document (Define-XML).

- Like cSDRG, ADRG also cannot have a summary of ADaM conformance findings, as the analysis datasets are not using CDISC ADaM format. The standard ADRG template can be used for filing.
- If R markdown coding is done for datasets and reports submission, then R Studio-related packages need to be mentioned with the version and functions used (project-specific/CRAN).
- Details about XPORT files and example codes for importing the XPORT files into SAS or R need to be written in the Data Format and Import Information section.

Define-XML (for tabulation and analysis datasets)

This file has a retrospective description of variable derivations. Challenges in creating RWD Define-XML: If there are any R datasets, they need to be translated into SAS datasets. To convert non-standard to CDISC-like format, adjust the Define stylesheet to accommodate longer dataset and variable names of legacy data for the same 'look and feel' of Define-XML v2.0.

Case Study

Have a pre-NDA meeting with the FDA

- Get clarification on how the FDA can perform analyses. They are usually open to work on SAS and R but would need deeper information on handling missing data or invalid data in the registry and procedure for variables derivation, which can be found in the SAP.
- The FDA asked for patient-level data to facilitate a complete review of the analysis results. And so data files were submitted in native format available, i.e. 'standard analytic file' metadata and data collection forms.
- They agreed to R markdown programs in submission but requested all R packages (versions) along with functions.
- The FDA requested data dictionaries.
- There was a proposal to submit v8 xpt because the datasets contain variables and labels of length greater than 8 and 40, respectively. But the FDA said it is not their policy to accept these file formats for xpts as version 8 has limitations, including increased file size, no native mechanism for support of audit trails, and referencing data sources. All electronic submissions for the NDA should use version 5 and follow the maximum permissible number of characters (8 characters for variable names, 40 for variable labels and 40 for dataset labels) based on the Study Data Technical Conformance Guide (July 2020, version 4.5.1).^{7,5} The FDA agreed to have both v8 and v5 xpts but with a linked document saying the dataset names, variable names and variable labels needed to be truncated, and mapping the original names to the shortened names.

- The eSUB package should be sent through an eCTD application so that the information is all in one place.

Submission package deliverables for RWD

- 1) Legacy tabulation datasets
 - aCRF
 - sas v8/v5 xpt files
 - Define files
 - Reviewer's guide
- 2) Legacy analysis datasets
 - sas v8/v5 xpt file – There was just one dataset with one record per transplant recipient, derived from intermediate datasets (.R files)
 - Define files
 - Reviewer's guide
- 3) Legacy programs
 - a) Dataset programs
 - sas to csv (.sas files) – Convert sas datasets to csv files
 - R markdown programs (.Rmd files) – Reads in csv files and creates intermediate R datasets (.Rdata files) and the final analysis dataset
 - R markdown reports (.html and .pdf files) – Results generated by programs
 - b) Output programs
 - Rmd programs that create R functions for analysis and tables/figures (PDF reports) in .Rmd files

7.6. Submission of RWD: What the Future Holds

Current FDA guidance^{7,2} dictates that clinical data must be submitted using the standards documented in the FDA Data Standards Catalog.^{7,6} This means, for the foreseeable future, CDISC will remain the de facto submission standard, even though it was custom-built for supporting RCTs and not quite fit for the diverse array of RWD. This requirement results in cumbersome data transformations, non-standard variables and domains, and business and validation rules that primarily apply to data collected in RCTs. The challenges of using CDISC standards for submitting RWD have been discussed by Jeff Abolafia et al in their papers.^{7,7,8,7,9}

As a viable alternative, the Observational Medical Outcomes Partnership – Common Data Model (OMOP CDM) or HL7 Fast Healthcare Interoperability Resources (FHIR) could be considered and developed to be submission-ready. There are many advantages to using either of the two. Not least, they both capture RWD and other non-interventional study designs optimally. FHIR is optimised for EHR and claims data, while the OMOP CDM is tailor-made for registry and other observational study designs.

However, there may need to be a paradigm shift in future regulatory submission standards wherein the argument should move away from a single submission standard to a more hybrid submission approach. This approach needs considerable effort from industry, regulators and standards authorities to harmonise and develop interoperability among data standards.

8. Glossary

- **RWD: Real-world data** is data relating to patient health status and/or the delivery of healthcare routinely collected from sources such as electronic health records, medical claims data, data from product or disease registries, and from sources (such as digital health technologies) that inform on health status.
- **RWE: Real-world evidence** is clinical evidence about the use and potential benefits or risks of a medical product derived from analysis of RWD.
- **Estimator:** In statistics, an estimator is a rule or formula to estimate an unknown quantity or parameter based on observed data. It's a method to make an educated guess about a population parameter, such as the mean or proportion, using sample data.
- **HTT:** Hypothetical target trial.
- **Length and Frequency of Follow-Up:** In research studies, especially longitudinal or observational studies, the length of follow-up refers to the duration over which participants are observed or tracked after the initial assessment or intervention. The frequency of follow-up indicates how often data collection occurs during that period. Both length and frequency of follow-up are crucial for understanding the trajectory of outcomes or changes over time.
- **Sample Size:** Sample size refers to the number of individuals or units included in a study or experiment. It's a fundamental aspect of study design, as it affects the reliability and generalisability of the results. A larger sample size generally provides more precise estimates and enhances the statistical power of the study to detect meaningful effects or differences between groups.
- **SPIFD:** Structured Process to Identify Fit-For-Purpose Data.
- **SPIFD2:** Structured Process to Identify Fit-For-Purpose Study Design and Data. Framework to Generate Valid and Transparent Real-World Evidence.
- **SPACE:** Structured Pre-Approval and Post-Approval Comparative Study Design.
- **StART-RWE:** Structured Template and Reporting Tool for Real-World Evidence.
- **Time Zero:** Time zero, also known as baseline, is the starting point or initial measurement in a study or experiment. It's the moment when observations or measurements begin, often used as a reference point for comparing changes or outcomes over time.
- **Treatment Group Assignment:** In experimental studies, particularly in clinical trials, treatment group assignment refers to the process of allocating participants to groups receiving different treatments or interventions. It's how researchers decide who receives the treatment being tested and who serves as the control group.

9. Disclaimer

The opinions expressed in this document are those of the authors and should not be construed to represent the opinions of PHUSE members, respective companies/organisations or regulators' views or policies. The content in this document should not be interpreted as a data standard and/or information required by regulatory authorities. All content are subject to change as newer guidelines being published by regulatory authorities.

10. Appendices

Appendix 1.1: Checklist

a) Details of Clinical Study

- Type of Clinical study per FDA Framework
 - Interventional Study
 - Observational Studies
- Clinical Study Purpose²

- Randomised clinical trial which uses RWD to capture clinical outcomes relating to safety or effectiveness
- Single arm trial which uses RWD as an external control arm
- Observational studies, such as observational cohort or case-control, which generate RWE for supporting an efficacy supplement
- Clinical trial or observational study which uses RWD/RWE to fulfil a post-marketing requirement (PMR)

b) Real-World Data Sources and Categorisation^{10.1}

- Electronic Health/Medical Records
 - Input data in HL7 Fast Healthcare Interoperability Resources (FHIR) platform
 - Input data not in HL7 but which follows proprietary standards
 - Input data does not follow any specific platform or data standards
 - Input data is SDTM-like
 - Not applicable
- Medical Claims Data
 - Proprietary efforts taken to transform and curate the data to CDISC
 - No efforts taken to transform and curate the data to CDISC
 - Input data is SDTM-like
 - Not applicable
- Product or Disease Registry
 - Registry follows NIH common data elements
 - OHDSI OMOP CDM model
 - Registry follows other proprietary data standards
 - Registry does not follow any data standards
 - Input data is SDTM-like
 - Not applicable
- Data Obtained from Digital Health Technologies
 - Input data is SDTM-like
 - Not applicable

- Other Data Sources (online health community, social media data, quality of life data collected from other platforms)
- Data sources follow proprietary data standards
 - Data sources do not follow any data standards
- c) Data Curation and Compliance Process**
- i) Based on the source of the data, has your organisation established a data curation process?
- Yes No
- ii) Are you familiar with the routine of a data migration plan to enable timely transfer of data from RWD sources?
- Yes No
- iii) Do you have a process for transforming unstructured source data into structured source data?
- Yes No
- iv) Is there a process in your organisation for harmonising the structured data across the system?
- Yes No
- v) Coding system mapping – Does the curation process involve dictionary term mapping?
(e.g. Systematised Nomenclature of Medicine Clinical Terms (SNOMED CT) to ICD-10-CM)
- Yes No
- d) Communication with External Stakeholders**
- i) Based on the inclusion of RWD in the clinical study, is your business function in communication with regulatory bodies?
- Yes No
- ii) For an oncology study, have you referred to the FDA QCARD initiative for communication with the US FDA?^{10.2}
- Yes No
- iii) Have you decided the timepoints in clinical trial progression for communicating with regulatory bodies?^{10.3}
- Yes No
- iv) Are you familiar with the contractual terms pertaining to data format and transfer between your organisation and the external vendor providing RWD?
- Yes No
- v) Are you aware of the platform and technology that the external RWD vendor uses to transfer RWD to your company system?
- Yes No

Appendix 1.2: Case Examples

Check-a-1: Type of Study

Option Selected by User – Observational Study

Implications and Decision Factors:

- 1) Refer to FDA guidance^{10.4} to support marketing application. Sponsor may be required to schedule a type-C meeting through the existing IND product.
- 2) Sponsor is required to do thorough documentation and annotation of programming codes pertaining to real-world data.
- 3) Further checks need to be evaluated by understanding health authority data submission requirements.

Check-b-1: RWD Sources and Categorisation

Option Selected by User – Electronic Health Records

Input data not in HL7 but follows proprietary standards

Implications and Decision Factors:

Involvement in and awareness of the following data operations processes by the statistical reporting group:

- 1) Data standards followed by the RWD vendor
- 2) Traceability and provenance assessment plan of the source data
- 3) Data migration plan, and use of the platform (either proprietary or vendor-specific)
- 4) Data curation plan
- 5) Planning for communication with health authorities regarding the nature and source of the data, and the data curation plan
- 6) Including these factors when planning data submission timelines

Check-c: Data Curation and Compliance Process

iii) Do you have a process for transforming unstructured source data into structured source data?

Option Selected by User: Yes

Implications:

- 1) The statistical reporting group needs to assess the applicability of existing data transformation processes for selected data sources for the clinical study.
- 2) Conduct a gap assessment to ensure the existing standard and platform can handle the existing data transformation processes.
- 3) Evaluate and plan for any risks in these processes and factor in that time when planning data submission timelines.

11. References

- 1.1 FDA. (2024). Real-World Data: Assessing Electronic Health Records and Medical Claims Data to Support Regulatory Decision-Making for Drug and Biological Products. Guidance for Industry.
- 1.2 https://pharmatimes.com/web_exclusives/Six_2020_predictions_for_real-world_evidence_1329688/
- 1.3 https://www.ema.europa.eu/en/documents/report/real-world-evidence-framework-support-eu-regulatory-decision-making-report-experience-gained-regulator-led-studies-september-2021-february-2023_en.pdf
- 1.4 <https://www.fda.gov/science-research/real-world-evidence/fda-use-real-world-evidence-regulatory-decision-making>
- 1.5 <https://www.ema.europa.eu/en/guideline-registry-based-studies-scientific-guideline#current-effective-version-section>
- 1.6 Yap, T.A. et al. (2022). Application of Real World Data to External Control Groups in Oncology Clinical Trial Drug Development. *Frontier in Oncology*, 11:695936. www.frontiersin.org
- 1.7 <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/considerations-design-and-conduct-externally-controlled-trials-drug-and-biological-products>
- 2.1, 2.2, 2.3 ICH. (2024). ICH M14 Guideline on general principles on plan, design and analysis of pharmacoepidemiological studies that utilize real-world data for safety assessment of medicines, May 2024.
- 2.4 'Tools to consider when engaging health authorities', Real World Data Solutions. TransCelerate <https://www.transceleratebiopharmainc.com/assets/real-world-data-solutions/>
- 4.1 Buneman, P., Khanna, S., & Tan, W-C. (2001). Why and where: A characterization of data provenance. *International Conference on Database Theory*, 316–330. https://doi.org/10.1007/3-540-44503-X_20
- 4.2 Simmhan, Y.L., Plale, B., & Gannon, D. (2005). A survey of data provenance in e-science. *ACM SIGMOD Record*, 34(3), 31–36. <https://doi.org/10.1145/1084805.1084812>
- 4.3 El Emam, K., & Arbuckle, L. (2013). *Anonymizing Health Data: Case Studies and Methods to Get You Started*. O'Reilly Media.
- 4.4 Malin, B., Loukides, G., Benitez, K., & Clayton, E. W. (2011). Identifiability in biobanks: models, measures, and mitigation strategies. *Human Genetics*, 130, 383–392. <https://doi.org/10.1007/s00439-011-1042-5>
- 4.5 Tene, O., & Polonetsky, J. (2013). Big data for all: Privacy and user control in the age of analytics. *Northwestern Journal of Technology and Intellectual Property*, 11(5), 239–273.
- 4.6 Floridi, L., & Taddeo, M. (2016). What is data ethics? *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2083), 20160360. <https://doi.org/10.1098/rsta.2016.0360>
- 4.7 European Union. (2016). **Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016** on the protection of natural persons with regard to the processing of personal data and on the free movement of such data (General Data Protection Regulation). *Official Journal of the European Union*, L119, 1–88. <https://eur-lex.europa.eu/eli/reg/2016/679/oj>
- 4.8 Government of Canada. (2000). **Personal Information Protection and Electronic Documents Act (S.C. 2000, c. 5)**. Retrieved from <https://laws-lois.justice.gc.ca/eng/acts/P-8.6/>
- 4.9 Health Canada. (2019). *Public Release of Clinical Information: Guidance Document*. Retrieved from <https://www.canada.ca/en/health-canada/services/drug-health-product-review-approval/profile-public-release-clinical-information-guidance/document.html>
- 4.10 U.S. Department of Health and Human Services. (2000). **Standards for Privacy of Individually Identifiable Health Information; Final Rule (HIPAA Privacy Rule)**. *Federal Register*, 65(250), 82462–82829. <https://www.hhs.gov/hipaa/for-professionals/privacy/laws-regulations/index.html>
- 4.11 European Medicines Agency. (2022). *Clinical Trials Regulation*. Retrieved from <https://www.ema.europa.eu/en/human-regulatory-overview/research-development/clinical-trials-human-medicines/clinical-trials-regulation>
- 4.12 European Medicines Agency. (2014). **European Medicines Agency policy on publication of clinical data for medicinal products for human use (Policy 0070)**. EMA/240810/2013. https://www.ema.europa.eu/en/documents/other/policy-70-european-medicines-agency-policy-publication-clinical-data-medicinal-products-human-use_en.pdf. **European Medicines Agency. (2010). European Medicines Agency policy on access to documents (Policy 0043)**. EMA/110196/2006. https://www.ema.europa.eu/en/documents/other/policy-43-european-medicines-agency-policy-access-documents_en.pdf
- 4.13 Chu, A. M. Y., Lam, B. S. Y., Tiwari, A., & So, M. K. P. (2019). An Empirical Study of Applying Statistical Disclosure Control Methods to Public Health Research. *International Journal of Environmental Research and Public Health*, 16(22), 4519. <https://doi.org/10.3390/ijerph16224519>

- 4.14 Abowd, J. M. (2018). The U.S. Census Bureau adopts differential privacy. *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2867–2867. <https://doi.org/10.1145/3219819.3226070>
- 4.15 Dwork, C., McSherry, F., Nissim, K., & Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. *Theory of Cryptography Conference*, 265–284. https://doi.org/10.1007/11681878_14
- 4.16 Dwork, C., & Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4), 211–487. <https://doi.org/10.1561/04000000042>
- 4.17 Machanavajjhala, A., Kifer, D., Gehrke, J., & Venkatasubramanian, M. (2007). I-Diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1), 3. <https://doi.org/10.1145/1217299.1217302>
- 4.18 Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Schulte Nordholt, E., Spicer, K., & de Wolf, P.-P. (2012). *Statistical Disclosure Control*. Wiley.
- 4.19 Willenborg, L., & de Waal, T. (2001). *Elements of Statistical Disclosure Control*. Springer. <https://link.springer.com/book/10.1007/978-1-4613-0121-9>
- 4.20 Templ, M., Meindl, B., Kowarik, A., & Alfons, A. (2022). *Statistical Disclosure Control in Practice: Protecting Confidentiality in Survey Microdata*. Springer.
- 4.21 Sweeney, L. (2002). k-Anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(5), 557–570. <https://doi.org/10.1142/S0218488502001648>
- 4.22 Ohm, P. (2010). Broken promises of privacy: Responding to the surprising failure of anonymization. *UCLA Law Review*, 57, 1701–1777.
- 4.23 Li, N., Li, T., & Venkatasubramanian, S. (2007). t-Closeness: Privacy beyond k-anonymity and I-diversity. *IEEE 23rd International Conference on Data Engineering*, 106–115. <https://doi.org/10.1109/ICDE.2007.367856>
- 4.24 Garfinkel, S. L. (2015). De-identification of personal information. NISTIR 8053, National Institute of Standards and Technology. <https://doi.org/10.6028/NIST.IR.8053>
- 6.1 ICH. (2024). ICH M14 guideline on general principles on plan, design and analysis of pharmacoepidemiological studies that utilize real-world data for safety assessment of medicines – scientific guideline, May 2024.
- 6.2 FDA. (2024). Real-world data: Assessing electronic health records and medical claims data to support regulatory decision-making for drug and biological products – guidance for industry, July 2024.
- 6.3 Gatto, N.M., Reynolds R.F., & Campbell, U.B. (2019). A Structured Preapproval and Postapproval Comparative Study Design Framework to Generate Valid and Transparent Real-World Evidence for Regulatory Decisions. *Clin Pharmacol Ther.*, 106(1), 103–115.
- 6.4 Gatto, N.M., Campbell, U.B., Rubinstein, E., et al. (2022). The Structured Process to Identify Fit-For-Purpose Data: A Data Feasibility Assessment Framework. *Clin Pharmacol Ther.*, 111(1), 122–134.
- 6.5 A Structured Process to Identify Fit for Purpose Study Design and Data to Generate Valid and Transparent Real World Evidence for Regulatory Uses - Gatto - 2023 - Clinical Pharmacology & Therapeutics - Wiley Online Library
- 7.1 FDA Guidance (2018). Framework for FDA's Real-World Evidence Program. December 2018. <https://www.fda.gov/media/120060/download>
- 7.2 FDA Guidance. (2023). Data Standards for Drug and Biological Product Submissions Containing Real-World Data: Guidance for Industry. December 2023. <https://www.fda.gov/media/153341/download>
- 7.3, 7.4 Abolafia, J., Ferko, S., & Holt, I. (2024). Future Clinical Data Submission Standards: CDISC, FHIR, OMOP, or Hybrid Model. Paper presented at the PHUSE US Connect 2024, Bethesda.
- 7.5 FDA Guidance (2020). Study Data Technical Conformance Guide. November 2020. <https://www.fda.gov/media/143550/download>
- 7.6 FDA Guidance. (2023). Data Standards Catalog. April 2023. <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/data-standards-catalog>
- 7.7 Abolafia, J., Ferko, S., & Holt, I. (2022). Submission Standards for RWD: The Good, the Bad and the Ugly. Paper presented at the PHUSE EU Connect 2022, Belfast. https://phuse.s3.eu-central-1.amazonaws.com/Archive/2022/Connect/EU/Belfast/PRE_RE09.pdf
- 7.8 Ferko, S., Holt, I., & Abolafia, J. (2023). Challenges and Considerations for Submitting Real World Data. Paper presented at the PHUSE US Connect 2023, Orlando. https://phuse.s3.eu-central-1.amazonaws.com/Archive/2023/Connect/US/Florida/PRE_RE05.pdf
- 7.9 Abolafia, J., Ferko, S., & Holt, I. (2023). Submission Standards for Real World Data: Gaps, Limitations and Recommendations. Paper presented at the PHUSE EU Connect 2023, Birmingham. https://phuse.s3.eu-central-1.amazonaws.com/Archive/2023/Connect/EU/Birmingham/PAP_RE03.pdf
- 10.1 Data Standards for Drug and Biological Product Submissions Containing Real-World Data - Guidance for Industry, Draft guidance by U.S. FDA. Released October 2021.

- 10.2 FDA. Oncology Quality, Characterization and Assessment of Real-World Data (QCARD) Initiative.
- 10.3 “Tools to consider when engaging health authorities”, Real World Data Solutions - TransCelerate <https://www.transceleratebiopharmainc.com/assets/real-world-data-solutions/>
- 10.4 FDA Guidance: Considerations for the Use of Real-World Data and Real-World Evidence to Support Regulatory Decision-Making for Drug and Biological Products, Final Guidance, August 2023.

12. Project Contact Information

- Dhruba Sikdar; Chi Zhang; Matt Baldwin; Elena Valkanova; Xingshu Zhu; Sherry Zhong; Yen Phan; Jingying Zhou; James Joseph; Priyadarshini Tunga; Srinath Yarasi; Sowmya Gabbula
- Email workinggroups@phuse.global

13. Acknowledgements

The authors thank the PHUSE RWE Working Group Leads. Without their support and guidance, this project would not have been possible.

The authors also thank the Integrated Data Analysis and Reporting (IDAR) function at Johnson & Johnson Innovative Medicine (JJIM). The structure and content of this paper are generously borrowed from the RWE playbook IDAR–JJIM developed for internal use. For that, the authors are highly indebted to the Johnson & Johnson team who developed the playbook.

Sincere thanks to the Submitting Real World Data Working Group project, led by Parag Shiralkar, for providing the submission checklist.

Lastly, the authors are grateful to their respective organisations for allowing them to work on this project.