# Technical Support Center: Harmonizing Disparate Real-world Data (RWD)

**Author**: BenchChem Technical Support Team. **Date**: December 2025

| Compound of Interest | | |
|---|---|---|
| Compound Name: | RW | |
| Cat. No.: | B13389108 | Get Quote |

Welcome to the technical support center for researchers, scientists, and drug development professionals. This resource provides troubleshooting guides and frequently asked questions (FAQs) to address common challenges encountered when harmonizing disparate real-world data (**RW**D) sources.

# Section 1: Semantic Harmonization & Data Standardization

This section addresses issues related to inconsistent coding, terminology, and units of measurement across different data sources.

# Frequently Asked Questions (FAQs)

Q1: My **RW**D sources (EHR, claims) use different codes for the same diagnosis (e.g., ICD-9, ICD-10, SNOMED CT). How can I standardize them for a unified analysis?

A1: This is a classic semantic harmonization challenge. The solution involves mapping local or varied codes to a single, standard terminology.

- Strategy:

  - Select a Target Standard: Choose a standard ontology appropriate for your research, such as SNOMED CT for clinical findings or LOINC for laboratory tests.

- Use Existing Crosswalks: Leverage established mapping resources and crosswalks (e.g., from the National Library of Medicine) to translate between code systems (like ICD-10-CM to SNOMED CT).

- Algorithmic & Manual Mapping: For local or non-standard codes, use natural language processing (NLP) and algorithmic searches to suggest potential matches to the standard terminology.[1] However, clinical expert review is critical to verify these matches and prevent misclassification.[1]

- Create a Reusable Mapping File: Document all mappings in a version-controlled file. This ensures reproducibility and transparency.[2]

Q2: I'm working with laboratory data from multiple international sites, and the units of measurement for the same test are different (e.g., mg/dL vs. mmol/L for glucose). What is the best practice for standardization?

A2: Standardizing units is crucial for accurate analysis. A systematic approach is required to prevent data loss and ensure clinical validity.[1]

- Strategy:

  - Profile the Data: Identify all unique tests and their corresponding units as they appear in the source data.

  - Define a Standard Unit: For each lab test, select a single, internationally recognized unit of measurement (e.g., using LOINC as a guide).

  - Verify and Convert: Use established clinical conversion factors to transform values. It is critical to have clinical experts review the test names, specimen types, and units to ensure the conversions are appropriate.[1]

  - Handle Non-Convertible Units: Document and quarantine records with units that cannot be reliably converted. A study on standardizing liver function tests found that this approach successfully converted the vast majority of records, with only 1.1% being excluded.[1]

Q3: We are trying to harmonize data from five different health data standards (e.g., HL7 FHIR, OMOP, CDISC). How do we manage the conceptual differences between data elements that

have similar names but different definitions?

A3: This requires moving beyond simple name matching to a concept-based harmonization approach. The goal is to map the underlying meaning of each data element.[2]

- Strategy:

  - Identify Concepts: For each topic (e.g., gender, vital status), identify the underlying concept represented by data elements across the different standards.[2]

  - Cluster Similar Concepts: Group the concepts that are semantically equivalent. For example, concepts representing biological sex might be clustered separately from those representing gender identity.[2]

  - Construct Mappings: Create explicit mappings between these concept clusters. This provides a more robust and context-aware harmonization than direct element-to-element mapping.[2]

  - Use a Common Data Model (CDM): A powerful approach is to map all source data to a CDM like the Observational Medical Outcomes Partnership (OMOP) CDM. This provides a standardized structure and terminology, facilitating large-scale, reproducible analyses.

# Troubleshooting Guide

| Problem | Possible Cause | Recommended Solution |
|---|---|---|
| High rate of mapping failure when using automated tools for lab codes. | Ambiguous or non-standard local test names; tool lacks context. | Implement a semi-automated approach. Use algorithms for initial matching, but ensure final validation is performed by clinical experts who understand the context of each data source.[1][3] |
| After merging datasets, patient counts for a specific condition are unexpectedly low. | Inconsistent diagnostic codes were not fully harmonized, leading to fragmented cohorts. | Re-run the code mapping process. Ensure all relevant codes (e.g., ICD-9, ICD-10, SNOMED) are mapped to a single target concept. Use a tool to explore the hierarchy of the target ontology to include parent concepts if necessary. |
| Data for a key lab value appears bimodal or has outliers after unit standardization. | An incorrect conversion factor was applied, or a subset of data was not converted. | Isolate the problematic data points and trace them back to the source. Verify the original units and the conversion factor applied. Implement data quality checks post-conversion to flag values outside of clinically plausible ranges. |

# Section 2: Data Linkage & Patient Identity

This section covers challenges related to accurately and securely linking patient data from different sources.

# Frequently Asked Questions (FAQs)

Q1: What is the best practice for linking a patient's clinical trial data with their **RW**D from EHRs and claims while maintaining privacy?

A1: The industry standard is privacy-preserving record linkage (PPRL) using tokenization.

- Strategy:

  - Obtain Patient Consent: It is a best practice to get patient consent for **RW**D linkage upfront, even if the data is de-identified.[4] Consent rates are often high (around 85%).[4]

  - Collect Personally Identifiable Information (PII): Securely collect a consistent set of PII (e.g., name, date of birth, address) from trial participants.[4]

  - Tokenization: Use a third-party service to convert the PII into an encrypted, irreversible token (e.g., a HealthVerity ID or HVID).[4][5] This token replaces the direct identifiers.

  - Link Across Datasets: The same tokenization process is applied to other **RW**D sources (EHR, claims). Records with matching tokens can then be linked without exposing the underlying PII.[4][5] This method allows for the creation of a longitudinal patient journey.[6]

Q2: We are trying to link records between a hospital's EHR and a trauma registry. What are the main linkage methodologies?

A2: There are two primary methods for record linkage:

- Deterministic Linkage: This method matches records based on an exact match of a set of unique identifiers (e.g., medical record number, social security number). It is straightfo**rw**ard but can fail if there are any errors or variations in the identifiers.[7]

- Probabilistic Linkage: This method is more flexible and powerful. It calculates a match probability score based on the agreement and disagreement of several identifiers (e.g., name, date of birth, zip code).[7] A threshold is set to determine which pairs are considered a match, a non-match, or require manual review. This approach is more resilient to minor data entry errors.

## Experimental Protocol: Probabilistic Data Linkage

- Data Preparation: Select linking variables (e.g., first name, last name, DOB, zip code) present in both datasets. Clean and standardize these variables (e.g., convert names to uppercase, format dates consistently).

- Blocking: Divide the datasets into smaller, manageable blocks based on a variable that is unlikely to have errors (e.g., the first letter of the last name or state of residence). This reduces the number of pairwise comparisons needed.

- Pairwise Comparison: Within each block, compare all possible pairs of records from the two datasets.

- Calculate Agreement Weights: For each linking variable, calculate agreement and disagreement weights (m- and u-probabilities) based on their estimated reliability and frequency.

- Compute Total Score: For each record pair, sum the weights to get a total linkage score.

- Set Thresholds: Define two thresholds: an upper threshold above which pairs are considered definite matches, and a lower threshold below which pairs are considered definite non-matches. Pairs with scores between the thresholds are sent for manual review.

- Evaluate Linkage Quality: Assess the linkage quality by calculating metrics such as sensitivity, specificity, and positive predictive value on a manually reviewed sample.

# Section 3: Data Quality and Completeness

This section addresses common problems with the intrinsic quality of **RW**D, such as missingness, errors, and inconsistencies.

# Frequently Asked Questions (FAQs)

Q1: My EHR dataset has a significant amount of missing data for a key variable (e.g., Body Mass Index). What are my options for handling this?

A1: The choice of method depends on the extent and pattern of missingness. Ignoring it can introduce significant bias.[8]

- Strategy:

  - Assess the Missingness: Determine if the data is missing completely at random (MCAR), at random (MAR), or not at random (MNAR). This will guide your strategy.

- Complete Case Analysis: The simplest approach is to analyze only the records with complete data. This is acceptable for small amounts of MCAR data but can lead to bias and loss of statistical power other**rw**ise.

- Single Imputation: Replace missing values with a single value, such as the mean, median, or mode. This is easy to implement but underestimates variance.

- Multiple Imputation: This is often the preferred method. It involves creating multiple complete datasets by imputing the missing values based on the distributions of the observed data. The analysis is performed on each dataset, and the results are pooled. This approach provides more accurate standard errors.

- Use of Unstructured Data: For some variables, the missing information may be present in unstructured clinical notes.[9] Natural Language Processing (NLP) techniques can be used to extract this information and fill in the gaps.[9]

Q2: How can we assess if our **RW**D sources are "fit-for-purpose" for a regulatory submission?

A2: The FDA emphasizes two core pillars for **RW**D quality: relevance and reliability.[10]

- Strategy:

  - Assess Relevance:

    - Clear Research Question: Start with a well-defined research question that the data must be able to answer.[10]

    - Variable Coverage: Ensure the dataset contains the necessary data elements (e.g., exposures, outcomes, covariates) with sufficient detail and follow-up time.

  - Assess Reliability:

    - Data Provenance: Document the origin of the data and how it was collected.

    - Data Quality Audit: Conduct a thorough audit of the data for completeness, accuracy, and timeliness.[10][11] Implement automated quality checks and maintain audit trails. [10]

Tech Support

- Standardization: Ensure data collection and formatting are standardized across sites to maintain consistency.[10]
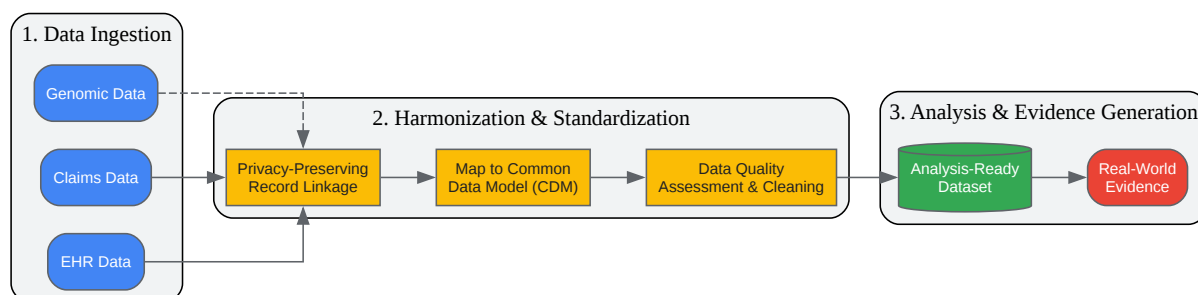
## Quantitative Data Summary

The following table summarizes common data quality issues and their potential business impact, highlighting the importance of addressing them early in the research lifecycle.

| Data Quality Issue | Description | Reported Business Impact | Common RWD Sources Affected |
|---|---|---|---|
| Inaccurate Data | Entries that are factually incorrect (e.g., misspelled names, wrong ZIP codes).[12] | The average organization loses an estimated $12.9 million annually due to poor data quality.[12][13] | EHR, Claims, Registries |
| Incomplete Data | Records with missing information in key fields (e.g., no value for BMI, missing race/ethnicity).[13][14] | Can lead to flawed analyses, unreliable conclusions, and biased results.[14][15] | EHR, Patient-Reported Outcomes |
| Duplicate Data | The same patient or event is recorded multiple times.[13][15] | Inflates patient counts, skews metrics, and increases storage costs.[13] | Claims, Registries |
| Data Heterogeneity | Data for the same concept is represented in different formats or codes.[16][17] | A significant barrier to integrating datasets and conducting multi-site studies.[8][16] | EHR, Lab Systems, Claims |

# Visualizations & Workflows
## General RWD Harmonization Workflow

The following diagram illustrates a typical workflow for harmonizing disparate **RW**D sources into an analysis-ready dataset.
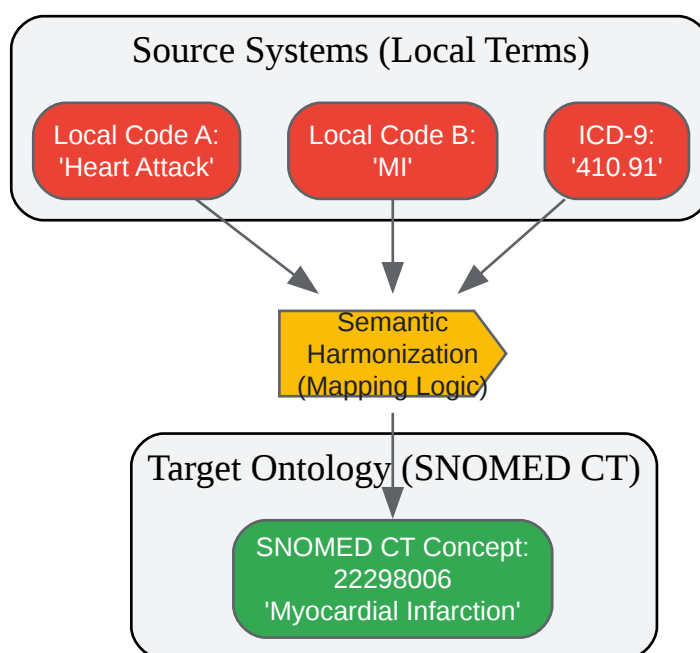


Click to download full resolution via product page

Caption: A workflow for harmonizing disparate Real-World Data sources.

## Semantic Mapping Logic

This diagram shows the logical relationship when mapping local, non-standard terms to a common ontology.



Click to download full resolution via product page

Caption: Mapping multiple local terms to a single standard ontology concept.

---

***Need Custom Synthesis?***

*BenchChem offers custom synthesis for rare earth carbides and specific isotopiclabeling.*

*Email: info@benchchem.com or Request Quote Online.*

---

# References

- 1. pharmasug.org [pharmasug.org]

- 2. researchgate.net [researchgate.net]

- 3. j2interactive.com [j2interactive.com]

- 4. blog.healthverity.com [blog.healthverity.com]

- 5. Linking clinical trial participants to their U.S. real-world data through tokenization: A practical guide - PMC [pmc.ncbi.nlm.nih.gov]

- 6. Unlocking Real-World Evidence Insights with Data Linking - Inovalon [inovalon.com]

- 7. Common Real-World Data Sources - Rethinking Clinical Trials [rethinkingclinicaltrials.org]

- 8. lifebit.ai [lifebit.ai]

- 9. ispor.org [ispor.org]

- 10. careevolution.com [careevolution.com]

- 11. 9 Common Data Quality Issues and How to Overcome Them [sagacitysolutions.co.uk]

- 12. firsteigen.com [firsteigen.com]

- 13. Gable Blog - 7 Common Data Quality Issues (and How to Solve Them) [gable.ai]

- 14. Real-world data: a comprehensive literature review on the barriers, challenges, and opportunities associated with their inclusion in the health technology assessment process - PMC [pmc.ncbi.nlm.nih.gov]

- 15. atlan.com [atlan.com]

- 16. Real World Data | CDISC [cdisc.org]

- 17. Breadth versus depth: balancing variables, sample size, and quality in Chinese cohort studies | The BMJ [bmj.com]

- To cite this document: BenchChem. [Technical Support Center: Harmonizing Disparate Real-world Data (RWD)]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b13389108#challenges-in-harmonizing-disparate-real-world-data-sources]

---

**Disclaimer & Data Validity:**

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

**Technical Support:** The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [Contact our Ph.D. Support Team for a compatibility check]

**Need Industrial/Bulk Grade?**   Request Custom Synthesis Quote

# BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: info@benchchem.com