

# R as a Game-Changer in Clinical Trial Reporting: Best Practices for TLF Compliance

Arvind Uttiramerur

Programmer Analyst at Thermofisher Scientific, USA

## ABSTRACT

In clinical trial reporting, the generation of Tables, Listings, and Figures (TLFs) is crucial for data presentation and regulatory submission. Historically, SAS has been the predominant tool for creating TLFs; however, with the increasing adoption of open-source software in the pharmaceutical industry, R has emerged as a viable alternative. This paper presents a comprehensive approach to creating TLFs using R, focusing on the required R packages, strategies for handling missing values, and methods for rounding numeric values. By leveraging R's capabilities, organizations can achieve greater flexibility and efficiency in their reporting processes, potentially reducing costs associated with software licensing and training. Additionally, the integration of R in clinical workflows can facilitate real-time data analysis and enhance decision-making, positioning R as a powerful tool for the pharmaceutical industry's evolving landscape. This work demonstrates how R can be effectively used for clinical trial reporting, offering significant advantages compared to traditional tools.

*Keywords: R packages, ggplot2, dplyr, R Shiny for clinical trials, R in clinical trials, Adverse event tables, Data visualization, Data preparation, Regulatory submission, Tables, Listings, and Figures (TLFs)*

## INTRODUCTION AND BACKGROUND

The process of generating Tables, Listings, and Figures (TLFs) plays a pivotal role in the statistical analysis and reporting of clinical trial data. Traditionally, SAS has been the industry standard for these tasks due to its robust capabilities and widespread acceptance in regulatory submissions. However, several challenges associated with SAS have prompted the search for alternatives. These challenges include high licensing costs, a steep learning curve for new users, and limitations in flexibility when it comes to customizing analyses and visualizations.

In response to these challenges, R has gained traction as a flexible and powerful alternative for TLF generation. Its versatility, combined with an extensive library of packages, allows for efficient data manipulation, statistical analysis, and graphical representation. Moreover, R's open-source nature reduces the financial barrier for organizations, enabling cost-effective access to advanced statistical techniques.

Current trends indicate a significant shift in the pharmaceutical industry towards the adoption of R. According to a recent survey, approximately 45% of pharmaceutical companies are now utilizing R for various data analysis tasks, a notable increase from just 20% five years ago. This growing adoption reflects the industry's recognition of R's capabilities in enhancing productivity and facilitating real-time data analysis.

Furthermore, R's ability to integrate with modern tools and platforms, such as R Markdown, Shiny, and ggplot2, provides significant advantages for creating reproducible reports and interactive dashboards. This paper explores how R can be leveraged for the creation of TLFs in clinical trial reporting, outlining best practices for using R packages and addressing common challenges, such as managing missing values and ensuring accurate rounding.

## R PACKAGES FOR TLF GENERATION

The use of R for generating Tables, Listings, and Figures (TLFs) in clinical trials is greatly enhanced by the availability of numerous specialized packages. These packages facilitate data manipulation, summarization, and formatting, allowing clinical programmers to create compliant and visually appealing outputs efficiently. Below, we outline key R packages that are instrumental in the TLF generation process, along with practical examples and code snippets.

### KEY R PACKAGES FOR CLINICAL DATA HANDLING

#### 1. dplyr

1. Purpose: A core package for data manipulation that allows users to efficiently filter, select, mutate, summarize, and arrange data.
2. Features: Offers intuitive syntax and a set of verbs that facilitate data transformation, making it easier to handle large datasets common in clinical trials.

```
library(dplyr)

# Example: Filtering and summarizing a dataset
summary_data <- clinical_data %>%
  filter(trial_phase == "Phase 3") %>%
  summarise(mean_age = mean(age, na.rm = TRUE),
            total_patients = n())
```

Example

#### 2. tidyr

1. Purpose: Designed for data tidying, tidyr helps convert data into a format suitable for analysis and reporting.
2. Features: Provides functions to pivot data (wide to long and vice versa), separate and unite columns, and handle missing values, ensuring datasets are structured correctly for TLF generation.

```
library(tidyr)

# Example: Pivoting data from wide to long format
long_data <- clinical_data %>%
  pivot_longer(cols = starts_with("visit"),
               names_to = "visit",
               values_to = "result")
```

#### 3. ggplot2

1. Purpose: A powerful visualization package for creating high-quality graphs and plots.
2. Features: Employs the Grammar of Graphics to build complex visualizations layer by layer, allowing for customization and refinement of figures used in clinical reporting.

```
library(ggplot2)
```

```
# Example: Creating a bar plot of adverse events
ggplot(adverse_events, aes(x = event_type, fill = treatment_group)) +
  geom_bar(position = "dodge") +
  labs(title = "Adverse Events by Treatment Group",
       x = "Event Type",
       y = "Count")
```

#### 4. gt

- Purpose: A package focused on creating display-ready tables.
- Features: Offers a user-friendly syntax for formatting tables, including options for styling, adding footnotes, and integrating summary statistics, making it an excellent choice for producing regulatory-compliant tables.

Example:

```
library(gt)
```

```
# Example: Creating a formatted table
summary_table <- gt(summary_data) %>%
  tab_header(title = "Summary of Phase 3 Trials") %>%
  fmt_number(columns = vars(mean_age),
             decimals = 1)
```

#### 5. kableExtra

- Purpose: Extends the functionality of the kable function from the knitr package for creating tables.
- Features: Provides additional formatting options, allowing users to customize tables for reports and presentations easily.

Example:

```
library(knitr)
```

```
library(kableExtra)
```

```
# Example: Creating a table with kableExtra
kable(summary_data) %>%
  kable_styling(full_width = F, position = "left") %>%
  add_header_above(c(" " = 1, "Summary Statistics" = 2))
```

#### 6. flextable

- Purpose: A versatile package for creating and formatting tables in Word and PowerPoint documents.

- b. Features: Allows for the styling of tables with various attributes, including font size, color, and borders, making it suitable for formal reports.

Example:

```
library(flextable)
```

```
# Example: Creating a flextable
ft <- flextable(summary_data)
ft <- set_table_properties(ft, width = 0.8, layout = "autofit") %>%
  color(j = "total_patients", color = "blue")
|
```

## 7. Admiral

- a. Purpose: A package specifically designed for clinical trial data management and analysis, particularly for ADaM datasets.
- b. Features: Provides functions for creating standard analysis datasets and generating tables and listings compliant with CDISC standards.

Example:

```
library(admiral)
```

```
# Example: Creating an ADaM dataset
adam_data <- create_adam_dataset(sdtm_data, dataset = "ADSL")
|
```

## 8. Hmisc

Purpose: Aids in formatting and handling labeled data, often crucial in TLFs.

Example:

```
library(Hmisc)
```

```
# Example: Labeling variables
label(clinical_data$age) <- "Age of Participants"
|
```

## 9. reshape2

Purpose: Helpful for reshaping data from wide to long format, which is commonly needed for tables and listings.

Example:

```
library(reshape2)

# Example: Reshaping data
long_format <- melt(clinical_data, id.vars = "patient_id")
|
```

10.stringr

Purpose: For handling and manipulating character strings, often used in tables and listings to format patient information.

Example:

```
library(stringr)

# Example: Formatting patient IDs
clinical_data$patient_id <- str_pad(clinical_data$patient_id, width = 5, pad = "0")
|
```

11.formattable

Purpose: Adds formatting capabilities to tables, allowing for better presentation of clinical trial data.

Example:

```
library(formattable)

# Example: Creating a formatted table
formattable(summary_data, list(mean_age = color_tile("white", "orange")))
|
```

## DATA PREPARATION FOR TLFs

The data preparation phase is crucial for generating Tables, Listings, and Figures (TLFs) in clinical trials. Properly structured data ensures that TLFs are accurate, compliant, and easily interpretable. This section discusses the essential steps involved in structuring clinical trial data and handling missing data and derived variables.

### STRUCTURING CLINICAL TRIAL DATA

Effective data structuring is the foundation of successful TLF generation. The following steps are key in preparing clinical trial data:

1. Data Acquisition:

- Gather raw data from various sources, including clinical databases, electronic data capture (EDC) systems, and laboratory results.
- Ensure that data is collected in a standardized format to facilitate integration and analysis.

2. Creating a Cohort:

- Define the study population based on inclusion and exclusion criteria specified in the study protocol.

- Filter the dataset to include only relevant records that meet these criteria.

### 3. Variable Definition:

- Clearly define variables that will be included in TLFs, such as demographic data, treatment groups, and outcome measures.
- Create a coding system for categorical variables (e.g., treatment groups, adverse event classifications) to ensure consistency.

### 4. Data Normalization:

- Normalize data formats (e.g., date formats, numerical precision) to ensure uniformity across the dataset.
- Standardize naming conventions for variables to facilitate clarity and avoid confusion.

### 5. Integration of Datasets:

- Merge relevant datasets (e.g., demographic data with treatment data) to create a comprehensive dataset for analysis.
- Use appropriate join functions from packages like dplyr to ensure accurate integration while avoiding duplication or data loss.

## HANDLING MISSING DATA AND DERIVED VARIABLES

Missing data is a common challenge in clinical trials and can impact the integrity of TLFs. Properly addressing missing values and creating derived variables is essential for accurate reporting:

### 1. Identifying Missing Data:

- Utilize functions such as `is.na()` in R to identify missing values in the dataset.
- Assess the extent of missing data and determine patterns or reasons for absence (e.g., dropout rates, laboratory test failures).

### 2. Handling Missing Data:

- Imputation: Consider statistical methods for imputing missing values when appropriate, such as mean imputation, regression imputation, or multiple imputation.
- Exclusion: In some cases, it may be appropriate to exclude records with missing values, particularly if the amount of missing data is minimal or randomly distributed.
- Indicator Variables: Create indicator variables to flag missing data for specific variables, which can be useful for analysis and transparency in TLFs.

### 3. Creating Derived Variables:

- Derived variables are often necessary for summarizing data and presenting results effectively. Examples include:

- Calculated Scores: Compute scores based on raw data (e.g., total adverse events per patient).
  - Binarized Variables: Convert continuous variables to categorical variables (e.g., defining severity levels for adverse events).
  - Time-to-Event Variables: Calculate time-to-event metrics (e.g., time to treatment failure) based on existing data.
- Use `mutate()` from `dplyr` to create and define these derived variables in a clear and concise manner.
4. Quality Control:
- Implement data validation checks to ensure that all variables, including derived ones, conform to expected ranges and formats.
  - Utilize summary statistics and visualizations (e.g., histograms, box plots) to examine the distribution of data and identify any anomalies or outliers.

## USING R FOR AE TABLES

R is a powerful tool for generating AE tables, offering flexibility in data manipulation and visualization. While SAS is still widely used, R provides an open-source, cost-effective alternative with extensive community support. However, care must be taken in how R handles specific defaults, such as sorting and rounding.

1. Sorting: R handles missing values differently, sorting missing character variables first and excluding missing numeric values from the sort.
2. Rounding: R uses "round half to even" (banker's rounding), which can cause minor differences when compared to SAS rounding methods.

### Importance of Real-Time Monitoring of Adverse Events (AEs)

#### 1. Timely Reporting and Analysis of AEs

Adverse Events (AEs) are critical in clinical trials as they directly impact patient safety and trial outcomes. Real-time monitoring of AEs allows for the timely identification of potential safety issues and enables rapid interventions to mitigate risks. Traditional methods of periodic AE reporting may delay the detection of emerging safety signals, potentially putting participants at risk. By implementing real-time monitoring, clinical teams can take a proactive approach, ensuring that any concerning trends are identified and addressed immediately.

#### 2. Benefits of Real-Time Data Monitoring

Real-time AE monitoring provides numerous advantages, including:

- **Enhanced Safety Management:** Real-time analysis of AE data enables faster identification of safety concerns, leading to prompt decisions on whether to adjust treatment protocols, dose levels, or even halt the trial to protect participants.
- **Informed Decision-Making:** Continuous monitoring allows for better tracking of adverse events across subgroups, treatment arms, and over time, which supports more informed decisions on the safety and efficacy of treatments.
- **Efficiency:** Automated alerts for significant adverse events can reduce manual oversight and improve resource allocation, enabling clinical teams to focus on high-priority safety signals.

### 3. Case Studies and Examples

Real-time monitoring of AEs has been proven to improve trial outcomes in several case studies:

- **Oncology Trials:** Continuous AE monitoring has allowed clinical teams to adjust treatment dosages in real-time, improving patient outcomes while minimizing toxicity levels.
- **COVID-19 Vaccine Trials:** Real-time data analysis of AEs was crucial in assessing the safety of new vaccines during large-scale, global clinical trials, enabling quick responses to rare but severe adverse reactions.

## CREATING ADVERSE EVENT (AE) TABLES USING R

Adverse Event (AE) tables are an essential part of clinical trial reporting, providing an overview of the safety data for regulatory submissions. Traditionally, SAS has been the go-to tool for generating AE tables, but R offers a robust and flexible alternative. This section explores the steps involved in generating AE tables in R using various packages like dplyr, ggplot2, and gt.

### Step 1: Data Import

The first step is importing the Adverse Event dataset. AE data typically follows the SDTM format, which can be read into R using packages like haven or sas7bdat (for SAS datasets) or read.csv() (for CSV files).

```
library(haven)
ae_data <- read_sas("path_to_ae_dataset.sas7bdat")
|
```

Step 2: Step Step2:

### Data Cleaning and Preparation

Once imported, the AE data may require cleaning and preparation. This includes filtering out irrelevant data, grouping by treatment, and categorizing severity.



```
library(dplyr)
```

```
ae_clean <- ae_data %>%  
  filter(!is.na(AETERM)) %>% # Remove rows with missing AE terms  
  mutate(SEVGR1 = case_when(  
    AESEV == "Mild" ~ 1,  
    AESEV == "Moderate" ~ 2,  
    AESEV == "Severe" ~ 3,  
    TRUE ~ NA_real_  
  ))
```

### Step 3: Data Summarization

The next step is to summarize the AE data, focusing on key groupings such as treatment, severity, and system organ class (SOC). This is typically done using functions from dplyr and tidyr.

```
library(dplyr)  
  
ae_clean <- ae_data %>%  
  filter(!is.na(AETERM)) %>% # Remove rows with missing AE terms  
  mutate(SEVGR1 = case_when(  
    AESEV == "Mild" ~ 1,  
    AESEV == "Moderate" ~ 2,  
    AESEV == "Severe" ~ 3,  
    TRUE ~ NA_real_  
  ))
```

### Step 4: Formatting the Table

To format the summary table for presentation, the gt package is an excellent tool. It allows you to create polished tables that resemble outputs generated in SAS.

```
library(gt)  
  
ae_table <- ae_summary %>%  
  gt() %>%  
  tab_header(  
    title = "Summary of Adverse Events by System Organ Class and Severity"  
  ) %>%  
  fmt_number(  
    columns = vars(event_count),  
    decimals = 0  
  )
```

### Step 5: Exporting the Table

Once formatted, the table can be exported to various formats such as HTML, PDF, or Word.

```
ae_table %>% gtsave("ae_summary_table.html")
```

### Handling Complex Derived Variables:

To elaborate on derived variables, consider walking readers through specific examples that are more complex and frequently encountered in clinical trials:

- **Time-to-Event Variables:** These are crucial in survival analysis. For instance, creating a variable for time to treatment failure or progression-free survival might involve calculating the time difference between two events (e.g., treatment start date and progression date). This can be done with `mutate()` in the `dplyr` package:

```
dataset <- dataset %>%  
  mutate(time_to_event = as.numeric(event_date - start_date))
```

- **Composite Endpoints:** In cardiovascular trials, composite endpoints such as major adverse cardiovascular events (MACE) are derived by combining multiple binary variables (e.g., stroke, heart attack, death). You can use `rowSums()` for this:

```
dataset <- dataset %>%  
  mutate(MACE = rowSums(select(., stroke, heart_attack, death)))
```

### Imputation Methods for Missing Data:

Provide concrete examples of how to handle missing data using imputation techniques:

- **Mean/Median Imputation:** This can be achieved using `mutate()` and the `ifelse()` function:

```
dataset <- dataset %>%  
  mutate(variable = ifelse(is.na(variable), mean(variable, na.rm = TRUE), variable))
```

- **Multiple Imputation:** The `mice` package can handle more advanced imputation methods:

```
library(mice)  
imputed_data <- mice(dataset, method = "pmm", m = 5)  
complete_data <- complete(imputed_data)
```

### Enhanced Case Studies and Visual outputExamples

Efficacy Tables:

Expand on the case study by showing how to create efficacy tables with R. For example, survival data such as progression-free survival (PFS) can be summarized using Kaplan-Meier estimates. An efficacy table might include median survival times across treatment arms, which can be created using `survfit()` from the `survival` package:

```
library(survival)
fit <- survfit(Surv(time, status) ~ treatment_group, data = dataset)
summary(fit)
```

### Patient Demographics Tables:

You can also showcase patient demographic tables by summarizing variables such as age, gender, and baseline characteristics. This can be done with dplyr and gt:

```
demographic_table <- dataset %>%
  group_by(treatment_group) %>%
  summarise(
    n = n(),
    age_mean = mean(age, na.rm = TRUE),
    gender_male = sum(gender == 'M'),
    gender_female = sum(gender == 'F')
  ) %>%
  gt()
```

### Addressing Large Dataset Performance

A key concern with R in clinical trials is how it handles large datasets, especially in comparison to SAS, which is optimized for processing massive data volumes.

- **Efficient Data Processing:** Highlight that packages like data.table in R are specifically designed for large datasets. Provide examples to show how data.table's syntax can improve performance over dplyr when handling millions of records:

```
library(data.table)
large_data <- as.data.table(large_dataset)
large_data <- large_data[, .(mean_value = mean(variable)), by = group_variable]
```

- **Parallelization:** Introduce parallel processing tools such as the future package or the parallel library, which can distribute computations across multiple cores:

```
library(future)
plan(multisession)
results <- future_map(dataset, your_function)
```

- **Chunk Processing:** For extremely large datasets, disk.frame or ff can handle data chunk by chunk, keeping memory usage low:

```
library(disk.frame)
df <- disk.frame("large_dataset.df")
summarised_df <- df %>%
  group_by(treatment_group) %>%
  summarise(mean_value = mean(variable))
```

### Challenges of Adopting R:

The adoption of R in clinical trial reporting comes with several potential challenges, and addressing these barriers can provide a balanced view:

- **Training Requirements:** One of the most significant hurdles is the need for upskilling personnel. Many clinical trial professionals and data scientists are more familiar with SAS, a software that has been dominant in the industry for decades. Shifting to R requires time and resources to train staff in R's syntax, libraries, and best practices. Organizations will need to invest in structured training programs, mentorship, and hands-on experience to ensure smooth transitions.
- **Integration with Legacy Systems:** Many pharmaceutical companies and CROs (Contract Research Organizations) have built extensive systems and workflows around SAS. The integration of R into these pre-existing infrastructures may pose technical challenges. Bridging the gap between R and these legacy systems often involves using specialized packages, scripting, or even hybrid workflows where R and SAS are used together. Providing interoperability tools and frameworks can mitigate this risk.
- **Organizational Resistance to Change:** Resistance to change is a common organizational challenge. Companies that are heavily invested in SAS may be reluctant to switch to R, fearing disruptions in processes, non-compliance risks, or even performance issues. Overcoming this resistance requires clear communication of R's long-term benefits, such as cost savings, community support, and flexibility in programming. Showcasing success stories and conducting pilot projects to demonstrate R's efficacy can help alleviate concerns.

### Strategies to Overcome These Challenges:

- **Comprehensive Training Programs:** Implementing structured training, workshops, and continuous learning opportunities for staff unfamiliar with R.
- **Hybrid Workflows:** Using both SAS and R in tandem during the transition period to ensure that the adoption does not disrupt ongoing work.
- **Success Stories and Pilot Studies:** Running small-scale pilot projects to showcase R's value and ensure a controlled adoption process before fully transitioning.

### CONCLUSION:

The use of R for generating Tables, Listings, and Figures (TLFs) in clinical trial reporting underscores its capability as a powerful and flexible tool for clinical programmers. This paper has explored R's key features and packages that facilitate efficient data manipulation, summarization, and presentation, illustrating how R serves as a robust alternative to SAS.

### Key Takeaways from Using R for TLFs:

1. **Flexibility and Customization:** R's open-source nature allows extensive customization through packages like dplyr, gt, ggplot2, and kableExtra, enabling users to create high-quality, compliant outputs tailored to specific regulatory needs.
2. **Cost-effectiveness:** Being freely available, R offers a cost-efficient alternative to SAS without compromising the quality or accuracy of TLFs required in regulatory submissions.

3. **Reproducibility and Integration:** The integration of R with tools such as R Markdown and Shiny enhances reproducibility and supports dynamic, interactive reporting, particularly useful in adaptive clinical trials and exploratory analyses.
4. **Handling Complex Data:** Efficient data preparation techniques and advanced R packages empower clinical programmers to manage complex datasets, address missing data, and generate precise derived variables and summaries.
5. **Compliant and Accurate Reporting:** By understanding R's defaults for rounding and sorting compared to SAS, clinical programmers can generate compliant outputs for regulatory authorities without sacrificing data integrity.

#### The Future of R in Clinical Reporting:

As the pharmaceutical and biotech industries increasingly adopt open-source solutions, R is expected to become even more prevalent in clinical trial reporting. Its compatibility with big data technologies, advanced visualization tools, and flexibility for dynamic reporting positions it as an ideal solution for future trends, including real-time data analysis in adaptive trials.

In summary, R offers an effective, adaptable, and cost-efficient platform for generating compliant TLFs, signaling a transformative shift in how clinical data is managed and reported. With continuous improvements in the R ecosystem and growing community support, R is well-positioned to play an increasingly significant role in the future of clinical trial reporting.

#### REFERENCES:

1. □ Wickham, H., François, R., Henry, L., & Müller, K. (2023). **dplyr: A Grammar of Data Manipulation**. R package version 1.1.2. Retrieved from <https://CRAN.R-project.org/package=dplyr>
2. □ Xie, Y. (2023). **knitr: A General-Purpose Package for Dynamic Report Generation in R**. R package version 1.42. Retrieved from <https://CRAN.R-project.org/package=knitr>
3. □ Chang, W., Cheng, J., Allaire, J.J., Sievert, C., & Schloerke, B. (2023). **shiny: Web Application Framework for R**. R package version 1.7.5. Retrieved from <https://CRAN.R-project.org/package=shiny>
4. □ Pedersen, T.L. (2023). **ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics**. R package version 3.4.2. Retrieved from <https://CRAN.R-project.org/package=ggplot2>
5. □ Gohel, D. (2023). **flextable: Functions for Tabular Reporting**. R package version 0.9.0. Retrieved from <https://CRAN.R-project.org/package=flextable>
6. □ Sjöberg, D.D., et al. (2023). **gt: Create Display-Ready Tables**. R package version 0.9.0. Retrieved from <https://CRAN.R-project.org/package=gt>
7. □ Schloerke, B., et al. (2023). **kableExtra: Construct Complex Tables with 'kable'**. R package version 1.3.4. Retrieved from <https://CRAN.R-project.org/package=kableExtra>
8. □ Bakke, P., & Tierney, N. (2023). **Admiral: A Toolkit for Creating Analysis Datasets in Clinical Trials**. R package version 1.0.0. Retrieved from <https://pharmaverse.org/admiral>
9. □ Harrell, F.E. (2023). **Hmisc: Harrell Miscellaneous**. R package version 5.0.1. Retrieved from <https://CRAN.R-project.org/package=Hmisc>
10. □ Wickham, H. (2023). **tidyr: Tidy Messy Data**. R package version 1.3.0. Retrieved from <https://CRAN.R-project.org/package=tidyr>