

# **R & Bioconductor in Drug Discovery**

SM

# Table of contents

<b>Preface</b>	<b>3</b>
<b>1 Introduction</b>	<b>4</b>
<b>2 Omics Technologies and Drug Discovery</b>	<b>6</b>
<b>3 R and Bioconductor in Omics: Comprehensive Support from Harmonization to Reporting</b>	<b>8</b>
3.1 1. Data Harmonization . . . . .	8
3.1.1 Annotation Packages . . . . .	8
3.1.2 Harmonization Tools . . . . .	8
3.2 2. Cross-Omics Analysis . . . . .	9
3.2.1 Statistical Frameworks . . . . .	9
3.2.2 Diverse Applications . . . . .	9
3.3 3. Tailored Visualizations . . . . .	9
3.3.1 Custom Visualization Tools . . . . .	9
3.3.2 Effect Size Representations . . . . .	10
3.4 4. Comprehensive Reporting . . . . .	10
3.4.1 Report Compilation . . . . .	10
3.5 5. Workflow Excellence . . . . .	10
3.6 Why R and Bioconductor? . . . . .	11
<b>4 Summary</b>	<b>12</b>
<b>References</b>	<b>13</b>

# Preface

Drug discovery is a complex and costly process, requiring innovative solutions to overcome challenges such as high expenses, long timelines, and the intricacies of biological systems. Computational tools like R and Python have become indispensable for addressing these hurdles, providing powerful methods to analyze and interpret vast amounts of biomedical data.

Both R and Python have increasingly adapted to one another, with Python tools incorporating R's philosophies and some R packages now having Python counterparts. This cross-adaptation enhances flexibility and enables seamless integration, allowing researchers to leverage the strengths of both languages in drug discovery.

In this guide, we'll explore how R and Bioconductor play a pivotal role in drug discovery. From single-cell RNA sequencing to multiomics and spatial transcriptomics, we'll showcase the tools and techniques used to analyze complex biological data. By examining real-world case studies and relevant packages, we aim to demonstrate how computational tools are driving innovation in pharmaceutical research and accelerating the development of new therapies.

# 1 Introduction

Drug discovery has evolved over centuries, from ancient remedies discovered by chance to modern research that identifies molecular targets linked to disease. Today, the process **begins with basic research identifying a therapeutic target (often a protein) associated with a disease. Researchers then search for therapeutic agents that modify the target to restore health.** However, finding the right target is challenging, as drugs often interact with unintended targets, sometimes with unknown effects.

The process progresses through a **preclinical phase (testing in animal models) followed by clinical trials in humans to establish drug safety and efficacy. If successful, regulatory agencies approve the drug, and it's released to the market. Post-approval, pharmacovigilance studies (phase 4) monitor long-term side effects.**

Drug discovery is a lengthy and expensive process—costing about **\$2.8 billion and taking 12–15 years** to complete. It generates vast amounts of data, often through high-throughput technologies, which can be analyzed with computational tools to speed up the process, reduce errors, and gain valuable insights.

This guide explores how **R and Bioconductor help analyze complex data from omics technologies—genomics, transcriptomics, proteomics, and metabolomics—enabling** researchers to uncover new drug targets and accelerate drug discovery as shown in this figure below from Singh et al. (2023).

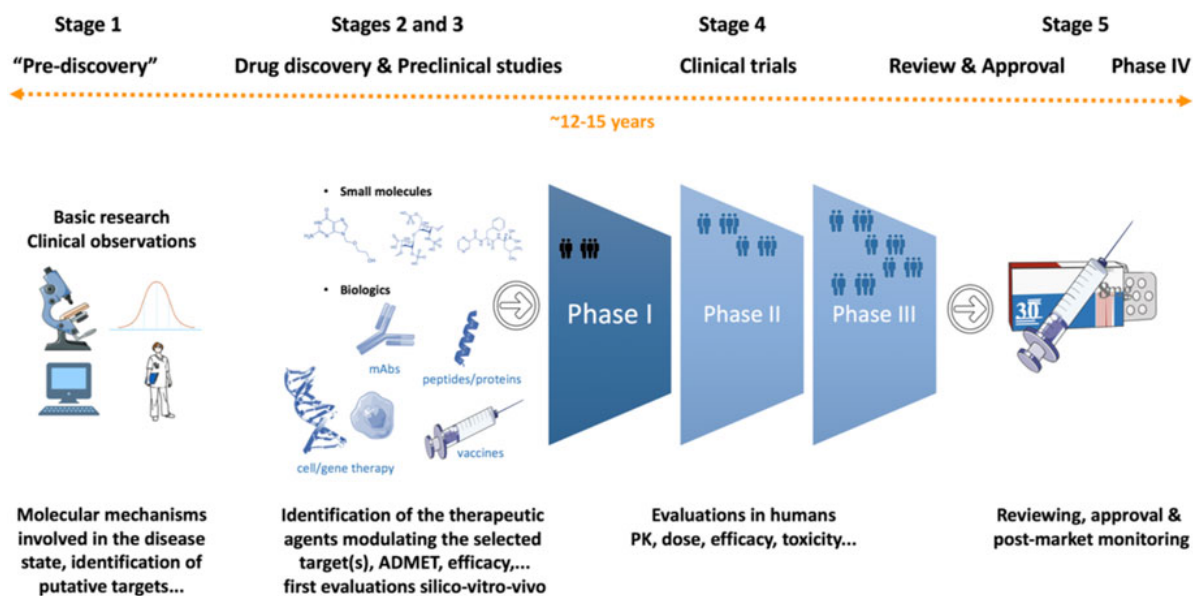


Figure 1.1: A simplified overview of the drug discovery and development process, which varies depending on the disease mechanisms and therapeutic agents

## 2 Omics Technologies and Drug Discovery

Omics technologies—such as genomics, transcriptomics, proteomics, and metabolomics—are revolutionizing drug discovery by enabling a more comprehensive and data-driven approach to understanding diseases. Here’s how each of these fields contributes to drug discovery

### 1. Genomics

**Role in Drug Discovery:** Genomics studies the entire genome, including genes, mutations, and regulatory elements. In drug discovery, genomics helps identify genetic variations associated with diseases, enabling the discovery of new drug targets (e.g., proteins or pathways altered in disease). **How Omics Help:** By sequencing genomes and identifying mutations, researchers can pinpoint genes responsible for diseases and design drugs that target these genes, potentially leading to precision medicine.

### 2. Transcriptomics

**Role in Drug Discovery:** Transcriptomics analyzes gene expression, i.e., which genes are turned on or off in a cell. This is crucial for identifying biomarkers and understanding how genes contribute to disease states. **How Omics Help:** By using technologies like RNA-sequencing (RNA-seq), researchers can observe the changes in gene expression patterns in response to drug treatments, helping identify the most promising drug candidates and biomarkers for disease monitoring.

### 3. Proteomics

**Role in Drug Discovery:** Proteomics focuses on proteins, their functions, and their interactions within a cell. Since proteins are the primary targets of most drugs, understanding the proteome can lead to better drug design. **How Omics Help:** Identifying protein markers or changes in protein function can help researchers design drugs that modify these proteins, improving treatment efficacy or targeting specific disease-related pathways.

### 4. Metabolomics

**Role in Drug Discovery:** Metabolomics studies the metabolites, which are the small molecules involved in metabolism. Since these metabolites reflect the biological activity of cells, they provide insights into cellular function and disease. **How Omics Help:** Analyzing metabolic changes can reveal how diseases alter metabolism and how drugs impact these pathways, facilitating the discovery of drugs that target metabolic disorders or influence disease progression.

## Integration of Omics in Drug Discovery

Omics technologies often work together to give a holistic view of a disease and its treatment. For example, combining genomics and transcriptomics can identify how genetic variations lead to changes in gene expression, which can be linked to the malfunctioning of certain proteins (proteomics). This comprehensive understanding allows researchers to develop drugs that can target multiple levels of biological processes.

Overview of omics data types and their functions in drug discovery, along with relevant databases supporting pharmaceutical research and development (adapted from @).

Omics	Function	Databases
<b>Genomics</b>	Understanding pathogenesis	GWAS Catalog
	Genetic association studies	GWAS central
	Identification of disease genes	dbGaP
	Discovery of putative drug targets	PharmGKB
	Patient-centered efficacy and toxicity assessment of drugs/targets	
<b>Transcriptomics</b>	Patient stratification	
	Disease mechanisms	DrugMatrix
	Mode of action of compounds	TG-GATE
	Moving from disease genes to drug targets	LINCS 1000
	Identification/evaluation of drug target candidates	Expression Atlas
<b>Proteomics</b>	Early prediction of adverse drug target effects	GEO repository
		ArrayExpress
	Post-translational process	PRIDE Archive
	Protein-protein network interaction	Peptide Atlas
	Drug target efficacy and safety evaluation at protein level	ProteomicsDB
<b>Metabolomics</b>		Human Proteome Map
	Protein toxicology	Human Proteome Atlas
	Novel DTD	Human Metabolome
	Drug target efficacy and safety evaluation at metabolomic level	Madison Metabolomics
	Metabolic toxicity	Golm Metabolome Database

## 3 R and Bioconductor in Omics: Comprehensive Support from Harmonization to Reporting

R and Bioconductor are indispensable for navigating the complexities of omics data, offering robust support at every step—from data harmonization to custom visualization and final report compilation.

---

### 3.1 1. Data Harmonization

Open-source databases like GWAS Catalog, GEO, and PRIDE Archive provide vast reservoirs of omics data. However, integrating data from these diverse sources requires harmonization to ensure consistency and compatibility.

#### 3.1.1 Annotation Packages

- **biomaRt**: Efficient retrieval of biomolecular annotations.
- **AnnotationHub**: Centralized resource for annotation datasets.
- **OrganismDbi**: Species-specific annotations.

#### 3.1.2 Harmonization Tools

- **GenomicRanges**: Aligning and harmonizing genomic coordinates.
  - **IRanges**: Handling interval data across datasets.
  - **BiocParallel**: Streamlining multi-core processing for harmonization tasks.
-



## 3.2 2. Cross-Omics Analysis

Once harmonized, the next step is analyzing the different omics data types to uncover meaningful insights.

### 3.2.1 Statistical Frameworks

- **Meta-analysis:**
  - **metaRNASeq:** Meta-analysis for RNA-seq studies.
  - **metafor:** General-purpose meta-analysis framework.
- **Integration:**
  - **mixOmics** and **MOFA2:** Multi-omics integration for holistic disease analysis.

### 3.2.2 Diverse Applications

- **Microbiome diversity** analysis with **phyloseq** to understand associations with phenotypes.
  - **Immune receptor sequencing** data interpretation using **Immunarch**.
- 

## 3.3 3. Tailored Visualizations

Visualizing complex results is critical for communicating findings effectively. R and Bioconductor provide unmatched tools for creating tailored, publication-ready graphics.

### 3.3.1 Custom Visualization Tools

- **ggplot2:** General-purpose visualizations.
- **ComplexHeatmap:** Multi-dimensional data representation.
- **Gviz:** Genomic data visualization.

### 3.3.2 Effect Size Representations

- Forest plots with **meta** and **metafor**.
  - Volcano plots and heatmaps for biomarker discovery.
- 

## 3.4 4. Comprehensive Reporting

R's capabilities extend beyond analysis and visualization to the creation of professional reports.

### 3.4.1 Report Compilation

- **R Markdown**: Seamless integration of code, results, and narratives.
  - **knitr** and **bookdown**: Automate reproducible workflows and produce high-quality books or documents.
  - **Quarto**: A modern framework for creating HTML, PDF, and other formats with enhanced styling.
- 

## 3.5 5. Workflow Excellence

Bioconductor provides comprehensive workflows for all major omics types, guiding users from raw, machine-generated output files to polished reports:

- **Genomics**: Starting with FASTQ/BAM files, workflows include alignment, variant calling, and visualization.
  - **Transcriptomics**: Covering raw RNA-seq data processing to differential expression analysis.
  - **Proteomics and Metabolomics**: Tailored pipelines for mass spectrometry data interpretation and pathway enrichment.
-

## 3.6 Why R and Bioconductor?

R and Bioconductor stand out as the only ecosystem offering an end-to-end solution for omics research: 1. **Harmonization and annotation** of open-source datasets. 2. **Advanced statistical methods** tailored for omics. 3. Tools for **microbiome**, **immune receptor**, and **multi-omics** analysis. 4. **Comprehensive and customizable visualization**. 5. **Reproducible report generation** with R Markdown and Quarto.

Whether you're identifying biomarkers, narrowing down drug targets, or compiling a professional report, R and Bioconductor empower researchers to achieve their goals seamlessly and efficiently.

## 4 Summary

In summary, this book has no content whatsoever.

## References

Singh, Natesh, Philippe Vayer, Shivalika Tanwar, Jean-Luc Poyet, Katya Tsaioun, and Bruno O Villoutreix. 2023. “Drug Discovery and Development: Introduction to the General Public and Patient Groups.” *Frontiers in Drug Discovery* 3: 1201419.