# Omics Data Science

## Concepts • Methods • Tools

2026 Edition

Authored by JIANGUO (JEFF) XIA

# Table of Contents

# Omics Data Science

*Concepts, Methods, and Tools*

## Copyright © 2026 by Jianguo Xia

• • • • •

### Permissions & Licensing

For permission requests or inquiries regarding academic use, please contact the author at jianguo.xia@mcgill.ca. For commercial training and software integration inquiries, please visit www.xialab.ai.

### Disclaimers

General Disclaimer. The information contained in this book is for educational and research purposes only. While the author has made every effort to ensure the accuracy of the information at the time of publication, the author assumes no responsibility for errors, omissions, or contrary interpretations of the subject matter herein.

Software & Code Disclaimer. The tools and data science workflows provided in this book are provided "as is" without warranty of any kind, express or implied. The author shall not be liable for any damages resulting from the use of the scripts or software tools described in this work.

### Trademarks & Affiliations

This work is published independently by the author. XiaLab Analytics, *MetaboAnalyst*, *ExpressAnalyst*, *MicrobiomeAnalyst*, *miRNet*, *OmicsNet*, *OmicsAnalyst*, and other associated software names and logos are trademarks or registered trademarks of XiaLab Analytics or their respective owners. The use of these names in this textbook is for educational and instructional purposes only and does not imply an endorsement of the book by the trademark owners, nor does it imply ownership of said trademarks by the author.

• • • • •

### Production Credits

First Edition: March 2026
Author: Jianguo Xia
Published in: Sainte-Anne-de-Bellevue, Quebec, Canada

ISBN: [Pending]
Cover Design: [Pending]
Interior Layout: [Pending]

• • • • •

### How to Cite This Book

This book was created to support the mission of **XiaLab Analytics** (https://www.xialab.ai): empowering researchers through training, tools, and AI. If you find this book helpful for your research, please cite it as:

*Xia, J. (2026). Omics Data Science: Concepts, Methods, and Tools.* Published by the Author. https://www.xialab.ca

# Preface

In June 2003, the Human Genome Project declared its work essentially complete. Thirteen years of international collaboration and roughly 2.7 billion dollars had produced a single reference genome. At the time, sequencing one human genome cost on the order of one hundred million dollars. Today, thanks to **next-generation sequencing (NGS)** technology, the same task costs well under a thousand dollars. Parallel revolutions in **mass spectrometry (MS)** have been equally dramatic. Comprehensive proteomic and metabolomic profiles that once required months of instrument time can now be generated in days. Any research group with a well-designed study can now produce multi-layered molecular portraits of health and disease.

Yet this extraordinary capacity to generate data has exposed a new bottleneck. The rate-limiting step in translating molecular measurements into biological understanding is no longer data generation. We now have unprecedented "parts lists" of genes, transcripts, proteins, metabolites, and microbes. But having a parts list is not the same as understanding how the system works. **The major bottleneck has shifted from data collection to data analysis.** Meeting that challenge requires a new kind of literacy: the ability to ask the right biological questions and to choose the right methods for understanding high-dimensional omics data.

Whether you are working with RNA-seq reads, proteomics, metabolomics, or multi-omic profiles, and whether your focus is biology, immunology, environment, or the microbiome, the analytical path remains remarkably consistent. The omics journey follows **a universal workflow**: converting raw instrument output into structured data, normalizing and filtering for quality, applying statistical methods to extract signal from noise, and finally, interpreting results within biological contexts. This book is a guide to that journey

## The Gap This Book Fills

Omics data analysis sits at an unusual crossroads. The omics technologies have become remarkably accessible: a single lab can now generate transcriptomic, proteomic, metabolomic, and microbiome data from the same set of samples in a matter of weeks. The bioinformatics tools have matured as well, with user-friendly platforms replacing what once required months of custom scripting. Yet a persistent gap remains between generating data and truly understanding it. This gap is not primarily technical. Most researchers can learn to click through a software pipeline. The real challenge is conceptual:

- Knowing **why** each data analysis step exists
- Understanding what **assumptions** it makes
- Recognizing **when** it fails

- Learning how to **choose** among different approaches

Without this level of data literacy, researchers risk generating polished but spurious results — mistaking statistical artifacts for genuine biological discovery. Yet no comprehensive **textbook** has existed to teach omics data science as an integrated discipline — one that weaves together the concepts, the statistics, and the practical know-how across all major omics types. This book was written to close that gap.

But conceptual understanding alone is not enough. A second, equally important gap has been the lack of easy-to-use, integrated bioinformatics **tools**. Too often, researchers who grasp the concepts find themselves unable to apply them — either because the right tools are scattered across dozens of packages that each require programming expertise, or because learning to code becomes a barrier that eclipses the science itself. Compounding this problem is a third challenge: the **sustainability** crisis in academic software. Most bioinformatics tools are published alongside a single paper and then gradually abandoned as their developers move on, leaving researchers stranded when the tools they depend on stop working.

The XiaLab tool ecosystem — including the freely available web-based platforms (*MetaboAnalyst*, *ExpressAnalyst*, *ProteoAnalyst*, *MicrobiomeAnalyst*, *NetworkAnalyst*, and others) and the comprehensive XiaLab Analytics Pro suite — was developed to address the three challenges. These tools provide intuitive, point-and-click interfaces where researchers can learn the methods described in this book and immediately apply them to their own data, without writing a single line of code. Equally important, they are maintained and evolved as a long-term commitment to the research community. The goal has always been to let scientists focus on the biology, not the software.

## Who This Book Is For

This book is for all life science researchers, regardless of computational background. It integrates the concepts of omics data science with accessible, hands-on tools. The primary audience is researchers and graduate students who work with omics data but do not necessarily have formal training in statistics or computer science.

You may be:

- A **biologist** who has just generated your first RNA-seq dataset
- A **clinical researcher** designing a multi-omics biomarker study
- A **bioinformatician** who is proficient in one omics type and wants to expand into others

You may prefer graphical interfaces, or you may write your own R scripts — either way, the analytical principles are the same.

We assume familiarity with basic biology (the central dogma, gene expression, metabolic pathways, the microbiome) and introductory statistics (means, $p$-values, hypothesis testing). Everything beyond that is developed from first principles within the text.

## How This Book Is Organized

The structure of this book mirrors the way omics analysis is actually done.

**Chapters 1–2** establish the foundation.

- **Chapter 1** introduces the conceptual framework that underlies all omics analysis: the two distinct challenges (size and complexity), the four-step analytical workflow — processing, normalization, statistical analysis, functional interpretation — the data distillation process that progressively reduces thousands of features to a handful of biological insights.

- **Chapter 2** introduces the core workflows for both sequencing (NGS) and mass spectrometry (MS): transforming raw instrument data into structured feature tables, and interpreting the resulting features within the context of molecular sets, pathways, and networks.

**Chapters 3–10** cover the individual omics disciplines. Most chapter pairs follow a consistent pattern: the first chapter introduces a discipline's fundamentals, and the second extends into advanced topics.

- **Chapters 3–4** cover transcriptomics — from RNA-seq reads mapping, differential expression to co-expression networks, dose–response analysis, and meta-analysis.

- **Chapter 5** introduces biological network analysis as both a conceptual framework and a practical tool for interpreting biological relationships.

- **Chapter 6** covers proteomics and biomarker analysis, with particular attention to the challenges of missing data and clinical validation.

- **Chapters 7–8** address metabolomics — from targeted quantification of known compounds to the untargeted discovery of unknown features.

- **Chapters 9–10** tackle the microbiome, covering amplicon-based community profiling, shotgun metagenomics, and the integration of microbial and metabolomics data.

**Chapters 11–12** bring everything together through multi-omics integration.

- **Chapter 11** details knowledge-driven approaches: focusing on how to contextualize omics findings within curated biological networks and molecular interaction maps.

**Chapter 12** presents data-driven approaches: dimensionality reduction, cross-omics correlation, clustering, and emerging methods including deep learning.

Together, these final two chapters provide the complete toolkit for combining multiple molecular layers into unified biological narratives.

## Guiding Principles

Three principles have guided the writing of every chapter:

- **Methods change; principles endure (Why before How).** For every analytical method, we explain the biological or statistical reason it exists before describing how to apply it. This matters more than knowing which button to click. Every method has assumptions, edge cases, and failure modes. We discuss these openly and share tips drawn from our decades of experience in data analysis.
- **Streamlined workflow patterns.** Omics analysis is a process of progressive refinement. Each chapter follows the natural flow from raw data (large, messy, technology-specific) to biological insight (compact, interpretable, biologically grounded). This pattern recurs throughout the book because it reflects how good analysis is actually done.
- **Hands-on practice.** To bridge theory and application, this book is designed to work in tandem with our companion slides and training videos, which demonstrate how to perform each analysis step-by-step using XiaLab tools.

## A Note on the XiaLab Tool Ecosystem

Throughout this book, we illustrate analytical workflows using the suite of web-based tools developed by the XiaLab:

- *ExpressAnalyst* for transcriptomics
- *ProteoAnalyst* for proteomics
- *MetaboAnalyst* for metabolomics
- *MicrobiomeAnalyst* for microbiome data
- *NetworkAnalyst* and *miRNet* for biological networks
- *OmicsNet* for knowledge-driven multi-omics integration
- *OmicsAnalyst* for data-driven multi-omics integration

These community tools are freely available and require no programming, making them accessible to researchers of all computational backgrounds.

That said, the analytical principles discussed in this book are tool-agnostic. Whether you use our platforms, R/*Bioconductor* packages, Python libraries, or commercial software, the same statistical foundations apply — and the knowledge you gain here transfers readily to any analytical environment.

## How to Use This Book

This book can be read in two ways:

- **As a course textbook:** Read the foundational Chapters 1 and 2 sequentially, as they build the core vocabulary and concepts. The remaining chapters cover distinct omics layers; they are designed to be independent yet complementary, allowing you to explore them in any order to fit a specific syllabus or build a complete systems-biology perspective.
- **As a reference:** Jump to the chapters you need. A metabolomics researcher, for example, can go directly to Chapters 7 and 8 after reading Chapters 1 and 2, without working through the transcriptomics or proteomics chapters first.

Each chapter follows a consistent format: learning objectives, narrative sections with callout boxes, figure placeholders, a summary, key terms, review questions, and further reading. The review questions are designed not for rote recall but for critical thinking — they ask you to evaluate scenarios, diagnose problems, and justify analytical choices.

Four types of callout boxes appear throughout:

- **[DEFINITION]** — Introduces a key term or concept at its first substantive appearance.
- **[TIP]** — Provides practical advice, shortcuts, or contextual insights.
- **[CAUTION]** — Warns about common mistakes, pitfalls, or misinterpretations.
- **[BEST PRACTICE]** — Recommends approaches and standard operating procedures.

## Acknowledgments

I thank the students in my bioinformatics courses, whose questions shaped how I think about teaching this material. The best explanations in this book are the ones that survived the most student questions; the remaining unclear passages are mine to improve in future editions.

I am grateful to colleagues who reviewed drafts and provided expert feedback across the diverse disciplines covered here. Any errors that remain are my responsibility.

Finally, I thank my family for their patience and support throughout the writing of this book.

<div align="right">

**Jianguo (Jeff) Xia, PhD**
*Professor, McGill University*
*Sainte-Anne-de-Bellevue, Quebec*
*March 2026*

</div>

# Chapter 1: The Omics Landscape

## Learning Objectives

By the end of this chapter, you will be able to:

- Explain the central dogma and how each omics layer — genomics, transcriptomics, proteomics, metabolomics — captures different biological information, including the microbiome as our second genome.
- Identify the two distinct challenges in omics data analysis (the size challenge and the complexity challenge) and explain why they require different solutions.
- Explain the role of rigorous study design — including biological replication, randomization, and blocking — in maximizing statistical power and controlling for bias across all omics experiments.
- Describe the data distillation process and the four general steps in omics data analytics.
- Distinguish between statistics, machine learning, and artificial intelligence in the context of sample size considerations and analytical goals.
- Understand key concepts including normalization, feature tables, batch effects, multiple testing correction, overfitting, and cross-validation.
- Navigate the XiaLab tool ecosystem and select appropriate tools for your data type.

• • • • •

## Introduction

The suffix "-omics" denotes the comprehensive study of a class of biological molecules. Just as genomics studies all genes, transcriptomics studies all transcripts, proteomics studies all proteins, and metabolomics studies all metabolites. Together, these approaches provide a multi-layered view of biological systems. They enable researchers to move beyond studying individual components and toward understanding how entire molecular systems function in health and disease.

> [DEFINITION: Omics] The comprehensive, high-throughput study of a complete set of biological molecules (genes, transcripts, proteins, metabolites, *etc.*) within a biological system.

The omics revolution has a specific origin and a breathtaking trajectory. When the Human Genome Project delivered its essentially complete reference sequence in 2003, sequencing a single human genome required thirteen years and roughly 2.7 billion dollars. Today, the same task takes hours and costs under a thousand dollars. Parallel revolutions have swept through mass spectrometry, enabling comprehensive proteomic and metabolomic profiling. What once required international consortia and dedicated genome centers can now be accomplished by individual laboratories in a matter of days.

Yet the ability to generate comprehensive molecular profiles is not the same as the ability to understand them. The Human Genome Project gave us a complete parts list of roughly 20,000 protein-coding genes, but a parts list does not explain how the system works — any more than an inventory of engine components explains how a car drives. The critical gap today is analytical: we must bridge the vast distance between raw molecular measurements and true biological understanding. Rigorous data science — the focus of this entire book — is how we build that bridge.

> **[TIP: From Parts Lists to Understanding]** The Human Genome Project provided the blueprint; subsequent omics technologies measure which parts are active, in what quantities, and under what conditions. The statistical and computational methods presented in this book are the tools that move us from cataloguing parts to understanding how those parts interact to produce health and disease.

This book is written for **busy researchers** who have data — often a great deal of it — and need a strategy for turning that data into biological insight. Perhaps you have received gigabytes of sequencing files from a genome center and are wondering how to get started. Perhaps you have a metabolomics feature table with over 10,000 peaks and 50 samples, and you need to find the handful of pathways that explain the difference between your treatment groups. Perhaps you are planning a multi-omics study and want to understand which analytical approaches apply to which data types.

Whatever your starting point, the journey follows the same fundamental pattern. Raw instrument data is progressively distilled — through quality control, data filtering, normalization, statistical analysis, and biological interpretation — until tens of thousands of measurements are reduced to a small number of actionable biological insights. This chapter introduces the conceptual framework for that journey. It establishes the core principles, vocabulary, and analytical strategies that every subsequent chapter builds upon.

# 1.1 The Extended Central Dogma

The classical central dogma describes information flow from DNA to RNA to protein. Modern omics extends this framework by capturing information at every layer of the molecular hierarchy. Each layer tells a different part of the biological story:

- **Genome: The blueprint.** DNA sequence is relatively static across the cells of an organism and changes slowly across generations. Genomics reveals what is **possible** — the full set of instructions encoded in the genome, including structural variants, copy number changes, and single-nucleotide polymorphisms (SNPs).

- **Microbiome: The Second Genome.** Most organisms do not exist in isolation but as part of a complex holobiont. The microbiome provides a vast, plastic reservoir of genetic diversity that complements the host. While the host genome is fixed, the microbial genome is dynamic—expanding metabolic capabilities, modulating immunity, and acting as a primary interface between the organism and its environment

- **Transcriptome: The active instructions.** Not all genes are active at any given time. Transcriptomics measures which genes are turned on or off, and at what levels. It provides a dynamic snapshot of cellular **intent** that changes with conditions, time, and cell type.

- **Proteome: The functional machinery.** Proteins are the workhorses of the cell, carrying out most biological functions. Proteomics captures what is **happening** by measuring protein abundance and post-translational modifications—such as phosphorylation or ubiquitination—that determine whether a protein is active or inactive.

- **Metabolome: The fuel and end products.** Metabolites are the dynamic small molecules reflecting the integrated output of all upstream layers and environmental factors. Metabolomics reveals what **actually happened**, representing the ultimate phenotypic consequence of the biological cascade

Each layer provides complementary information. A gene can be present but not expressed; a transcript can be abundant, but its protein product may be post-translationally modified and inactive; a metabolite reflects the integrated function of multiple pathways operating simultaneously. By studying these layers in concert, we gain a comprehensive understanding of biological systems that no single layer can provide.

The following diagram illustrates this hierarchical flow from genome to phenotype, with each layer becoming progressively more dynamic and each omics technology capturing a distinct slice of cellular activity.

*Figure 1.1 — The omics hierarchy from genome to phenotype. A vertical flow diagram showing: Genome/Microbiome (DNA) → Transcriptome (RNA) → Proteome (Proteins) → Metabolome (Metabolites) → Phenotype. Each layer is labeled with its characteristics: static to dynamic, blueprint to end products. Arrows indicate information flow and the omics technology that measures each layer.*

> **[TIP: Beyond the Linear Flow]** The central dogma is a functional simplification. In reality, biological systems are governed by complex feedback loops: metabolites provide the substrates for epigenetic modifications that regulate gene expression, proteins act as transcription factors to throttle their own production, and environmental exposures influence every layer simultaneously. Modern omics is increasingly focused on capturing these bidirectional, non-linear interactions

## 1.2 Two Distinct Challenges

Understanding omics data science begins with recognizing two fundamentally different challenges. These challenges require different approaches and different ways of thinking. Conflating them is one of the most common sources of confusion for newcomers to the field.

### Challenge 1: The Size Challenge (Raw Data)

The first challenge is sheer data volume. Omics instruments generate massive amounts of raw data:

- A single RNA-seq sample produces a FASTQ file of 1–10 GB.
- A high-resolution LC-MS metabolomics sample generates 100s MB of spectral data.

A study with 100 samples easily produces 100–500 GB of raw data, and large-scale projects routinely accumulate terabytes (TB). This volume creates a practical problem: researchers who have spent thousands of dollars on sequencing or MS receive huge datasets and face an immediate question — "How do I get the most from this data? How do I even get started?"

The good news is that raw data processing is technology-specific and largely solved. Each instrument type produces data in its own format, but standardized bioinformatics pipelines exist to convert these vendor formats into structured feature tables. Quality control metrics are well established, alignment algorithms are mature, and peak-picking software is increasingly reliable. The computational infrastructure required is still significant, but the intellectual challenge has been substantially addressed by the bioinformatics community over the past decades.

## Challenge 2: The Complexity Challenge (Feature Tables)

After raw data preprocessing, we obtain a **feature table** — a structured matrix where rows represent features (genes, proteins, metabolites, microbial taxa) and columns represent samples. Cell values are abundance measurements. This table is orders of magnitude smaller than the raw data, but it is inherently more complex to analyze.

> **[DEFINITION: Feature Table]** A data matrix where rows represent molecular features (genes, transcripts, proteins, metabolites, microbial taxa, or peaks) and columns represent samples. Cell values represent abundance, expression, or quantification levels. This is the universal starting point for omics data analysis.

The feature table appears structurally simple, yet it is profoundly challenging to analyze. Several intrinsic properties conspire to make omics data science fundamentally different from traditional experimental statistics:

- **High dimensionality:** A typical feature table profiles thousands to tens of thousands of distinct molecular features from each sample in a single experiment.

- **Small sample sizes:** Most studies have only tens to low hundreds of samples. This creates the "$p \gg n$" problem — more features than samples — which violates the assumptions of many traditional statistical methods.

- **Technical variation:** Batch effects from different processing dates, instruments, or reagent lots introduce systematic non-biological variation.

- **Biological variation:** Natural heterogeneity between individuals means the same condition can produce different molecular profiles in different subjects.

- **Missing values:** Not all features are detected in all samples, whether due to detection limits, dropout, or stochastic variation in measurement.

- **Non-normal distributions:** Abundance values are typically skewed, with many low-abundance features and a few highly abundant ones. This requires mathematical transformation before analysis.

Unlike raw data preprocessing (which is largely standardized), feature table analysis is open-ended. It supports diverse hypotheses, iterative exploration, and creative interpretation over extended periods. This is where the real science happens — and it is fundamentally more challenging than raw data processing.

> **[TIP: Why This Distinction Matters]** Recognizing these two challenges helps you understand why raw data processing can be automated and standardized, while feature table analysis requires scientific judgment, multiple methods, and biological validation. Chapter 2 covers the technology-specific solutions to the size challenge in detail.

The contrast between these two challenges is stark: raw data processing funnels gigabytes of instrument output into a compact feature table (a largely solved engineering problem), while feature table analysis confronts the high-dimensional, small-sample properties of that table head-on (an open-ended scientific problem).

*Figure 1.2 — The two challenges side by side. Left panel: Size Challenge showing raw data volume (FASTQ files, raw spectra) funneling into a compact feature table. Right panel: Complexity Challenge showing the feature table with its properties (high dimensionality, small n, missing values, batch effects) leading to analysis and biological interpretation. An arrow connects them: raw data processing (solved) feeds into feature table analysis (open-ended).*

## 1.3 Rigorous Study Design

The most sophisticated computational analysis cannot salvage a fundamentally flawed experimental design. Before a single sample is processed, the integrity of the downstream analysis is determined by three core principles: **biological replication**, **randomization**, and **blocking**. Applying these principles across all omics layers is the best defense against spurious findings and the most effective way to maximize statistical power.

### Biological Replication (Sample Size)

A **biological replicate** is an independent measurement from a distinct biological entity (e.g., a separate mouse, a different patient, an independent cell culture batch).

- **Why it matters:** Biological replication is necessary to measure the natural variation within a population. Without it, you cannot determine whether the difference between your control and treatment groups is a true biological effect or simply random individual variation. The required number of replicates depends on the expected effect size and the inherent variability of the biological system. Larger effect sizes and lower variability require fewer samples; small effects in noisy systems require many.

- **A common mistake:** Technical replicates (re-running the same sample multiple times) measure instrument precision but cannot substitute for true biological replicates.

## Randomization (Controlling for Bias)

**Randomization** is the process of assigning samples to treatment groups and processing order randomly.

- **Why it matters:** Randomization ensures that any unknown or unmeasured factors — such as minor environmental differences, the time of day a sample was collected, or a subtle trend in instrument performance — are evenly distributed across all experimental groups. This prevents bias from becoming confounded with your biological variable of interest (e.g., "treatment").

- **Example:** If the first 10 samples processed are all controls and the last 10 are all treatments, any subtle drift in the instrument over time becomes confounded with the treatment effect. Randomizing the processing order breaks this link.

## Blocking (Managing Known Variation)

**Blocking** is a design technique used to group samples based on a known source of technical or biological variation that is *not* the factor of interest.

- **Why it matters:** Blocking is essential for managing batch effects (as discussed in Section 1.7). If a study must be processed over two days, the "Day" variable is the block. By ensuring that each block contains a balanced set of all treatment groups, the statistical model can estimate and remove the systematic variation caused by the "Day" factor — isolating the pure biological effect of the treatment.

- **How to implement it:** A well-designed block balances the treatment groups across the technical factor. For instance, half of the control and half of the treated samples are processed on Day 1, and the remaining halves are processed on Day 2.

    > **[BEST PRACTICE: Design Before Experiment]** The time spent on rigorous study design — ensuring adequate biological replication, randomizing sample processing order, and explicitly blocking for known technical factors — is the most statistically powerful investment you can make in any omics experiment. No amount of downstream analysis can correct for a poorly designed study.

## 1.4 Omics Data Distillation

One of the most important conceptual frameworks in omics data science is the **data distillation process** — progressively filtering and condensing high-dimensional data into actionable biological insights. Think of it as a funnel: at the top, you pour in everything the instrument measured; at the bottom, you extract the essential biological story.

> **[DEFINITION: Omics Data Distillation]** The process of transforming large, high-dimensional datasets into a manageable set of biologically meaningful features through filtering, statistical analysis, and integration with biological knowledge.

The distillation follows a characteristic pattern, with the number of features decreasing at each stage:

**Raw Features: ~100,000** (sequencing reads, MS peaks)

↓

**Data Cleaning: ~10,000** (remove low-quality, uninformative features)

↓

**Statistical Filtering: ~1,000** (significant after multiple testing correction)

↓

**Pathway/Network Analysis: ~100** (aggregated into functional categories)

↓

**Biological Candidates: ~10** (actionable biomarkers or pathways)

↓

**Final Insights: < 5** (novel hypotheses and mechanisms)

A concrete example from metabolomics makes this tangible. In a study using high-resolution Orbitrap mass spectrometry, the raw data contained 35,000 features (ions and metabolites). After data cleaning and quality filtering, this was reduced to approximately 10,000 features. Statistical analysis identified around 500 significant features. Functional annotation and pathway analysis consolidated these into 15 significant metabolic pathways. Network analysis then revealed 2 central metabolic hubs, leading to mechanistic hypotheses and therapeutic targets.

The key insight is that each step in the funnel serves a specific purpose: quality filtering removes noise, statistical testing separates signal from background, and functional analysis translates molecular-level findings into biology. Skipping steps — or applying them in the wrong order — produces suboptimal or unreliable results.

The following diagram visualizes this progressive narrowing from raw measurements to final biological insight, with the methods used at each stage labeled alongside.

*Figure 1.3 — The omics data distillation funnel. A funnel-shaped diagram narrowing from top to bottom. At each level: raw features (100K), cleaned features (10K), statistically significant features (1K), pathway-annotated features (100), actionable candidates (10), and final insights (<5). The left side labels the process at each step; the right side labels the methods used (QC, statistics, enrichment, network analysis).*

## Three Strategies for Data Distillation

The omics community has evolved three complementary approaches for identifying signals in high-dimensional data. Understanding the differences between them is essential for choosing the right analytical strategy.

**1. Unsupervised exploration (data-driven).** This approach does not use outcome labels or metadata. It identifies intrinsic patterns through data exploration alone — using methods like **principal component analysis (PCA)** or clustering, coupled with visualization. Unsupervised methods are ideal for hypothesis generation and for understanding data structure before any formal testing. For example, you might discover that your samples cluster by disease status in a PCA plot before you even run a statistical test. Because these methods do not use outcome information, they carry no risk of overfitting to outcome labels.

**2. Supervised analysis (model-based).** This approach relies on outcome labels (e.g., disease vs. control) to drive feature selection and train predictive models. While methods ranging from classical statistics (*t*-tests, ANOVA) to machine learning classifiers are highly effective at identifying condition-specific features, they carry an inherent risk. In high-dimensional datasets with thousands of features but few samples, the model has enormous flexibility to find spurious

correlations that appear significant by pure chance. Consequently, rigorous **multiple testing correction** and **cross-validation** are absolutely essential.

**3. Contextualized analysis (knowledge-guided).** This approach combines the first two with external biological knowledge. After identifying statistically significant features, you ask whether they make biological sense: Are they known to participate in relevant pathways? Do they interact within known biological networks? Are they consistent with previous literature? This integration of data-driven results with biological context is the most effective strategy for distinguishing true signals from statistical artifacts.

> **[BEST PRACTICE: Start Unsupervised, Then Supervised, Then Contextualized]** In practice, the most effective approach applies these three strategies sequentially. Start with unsupervised exploration to understand your overall data structure and catch hidden issues like batch effects or outliers. Then, apply supervised methods to identify statistically significant features. Finally, validate those results by contextualizing them within known biological pathways. This layered approach is the hallmark of rigorous omics analysis.

## 1.5 Four General Steps in Omics Data Analytics

Despite the vast diversity of omics technologies—from next-generation sequencing to high-resolution mass spectrometry—nearly all workflows follow a remarkably consistent four-step framework. Understanding this "Universal Pipeline" provides a mental map for every analysis you will encounter in this book.

### Step 1: Raw Data Preprocessing (Technology-Specific)

The first step transforms raw instrument output into a structured feature table. This step is inherently technology-specific: sequencing data requires read alignment or assembly, while mass spectrometry data requires peak picking and spectral processing. The input is a vendor-specific file, and the output is a structured feature table with features in rows, samples in columns, and abundance values in cells.

This step is now largely standardized through established bioinformatics pipelines. Chapter 2 covers raw data preprocessing in detail for the two dominant technologies: next-generation sequencing (NGS) and liquid chromatography–mass spectrometry (LC-MS).

## Step 2: Data Processing and Normalization (Omics-Specific)

Once you have a feature table, it must be prepared for statistical analysis by addressing technical artifacts. This step is where many analyses go wrong, because technical variation can easily be mistaken for biological signal if not properly handled.

The key processes at this stage include:

- **Quality control:** Remove outlier samples, assess batch effects, and identify sample mix-ups. Visual inspection (PCA plots, box plots, density plots) is critical here.
- **Filtering:** Remove low-quality or non-informative features that add noise without contributing signal.
- **Missing value handling:** Features not detected in some samples must be addressed through imputation (replacement with estimated values) or removal.
- **Normalization**, which has two components:
    - **Technical normalization** adjusts for technical variation — differences in sequencing depth, sample amount, or instrument sensitivity — to ensure features are comparable across samples.
    - **Statistical normalization** applies mathematical transformations (such as log transformation or scaling) to address skewed distributions and stabilize variance.

Why does normalization matter? Consider a simple example. If one RNA-seq sample was sequenced to 40 million reads and another to 20 million reads, every gene in the first sample will appear to have roughly twice the expression of the same gene in the second sample. Without normalization, this sequencing depth difference creates a systematic false signal across the entire dataset. Normalization removes this technical artifact so that only true biological differences remain.

A few specialized methods (such as *limma* and *DESeq2*) provide their own built-in normalization optimized for statistical analysis. For most other statistical and machine learning methods, however, data that has been transformed to approximate a normal distribution tends to produce the best results. The appropriate choice depends on your data type and analytical context.

The following comparison illustrates how normalization aligns sample distributions, removing systematic technical differences while preserving true biological variation.

*Figure 1.4 — The effect of normalization. Two panels: Left shows box plots of sample distributions before normalization, with clear systematic differences in median and spread between samples processed in different batches. Right shows the same samples after normalization, with distributions aligned and comparable across samples.*

## Step 3: Statistical Analysis and Visualization (Omics-Generic)

With a properly normalized feature table in hand, the next step is to identify patterns and significant features using statistical and machine learning methods. This step is largely omics-generic — the same methods work whether your features are genes, proteins, or metabolites.

The analytical toolkit includes:

- **Univariate analysis:** This stage involves testing individual features for their association with specific outcomes. While classical methods like *t*-tests and ANOVA provide the foundation, we recommend modern omics-specific frameworks like *limma*, *edgeR*, or *DESeq2*. These tools are explicitly engineered to handle the unique statistical challenges of high-dimensional data, namely small sample sizes and unstable variance. By yielding both a formal *p*-value and an effect size (fold change), this step provides a clear picture of both the statistical certainty and the practical biological magnitude of every signal.

  - **Multiple testing correction:** When testing thousands of features simultaneously, some will appear significant by chance alone. The **false discovery rate (FDR)** correction, typically using the Benjamini–Hochberg method, controls the expected proportion of false positives among all features declared significant.

  - **Information borrowing:** With only a few replicates per group, the variance estimate for any single feature is unreliable. Moderated statistics most commonly – most commonly **empirical Bayes shrinkage** – solve this by "borrowing information" across the entire dataset. By adjusting extreme, feature-specific variances toward a global trend, these methods produce highly stable estimates even with very small sample sizes. This stabilization is exactly why specialized omics tools drastically outperform generic, feature-by-feature *t*-tests. Chapters 3 and 6 cover this mechanism in depth.

    > **[DEFINITION: *p*-value]** The probability of observing a test statistic as extreme as — or more extreme than — the observed value, assuming the null hypothesis is true (i.e., no real effect). A *p*-value describes the likelihood of random chance, **not** the strength of a biological signal. A *p*-value of 0.001 with a fold change of 1.05 indicates a tiny effect that, while statistically detectable, may not be biologically relevant.

- **Multivariate analysis:** This stage considers all features simultaneously to capture the complex biological relationships and interactions within the data. While univariate methods evaluate one feature at a time, multivariate techniques reveal the global

structure of the dataset that individual tests inevitably miss. We typically group these methods into three distinct categories.

- o **Dimensionality Reduction (Unsupervised)**: Before grouping or predicting, we must often simplify the data. Methods like Principal Component Analysis (PCA) condense thousands of features into a few major "principal components". This allows us to visualize the major axes of variation, assess overall data quality, and spot broad trends in a 2D or 3D plot without relying on outcome labels

> **[DEFINITION: Principal Component Analysis (PCA)]** PCA is the primary tool for simplifying high-dimensional omics data by compressing thousands of features into a few **Principal Components (PCs)**. Each PC is a **linear combination** of the original variables, meaning every feature is assigned a weight, or "loading", based on its contribution to that specific axis. This mathematical relationship allows for two-way interpretation: **Scores** represent the coordinates of samples in the new space (Figure 1.5, panel 1), while **Loadings** reveal which specific genes or proteins are driving the patterns seen in those samples. Because it is an unsupervised method, it is the ideal first step for identifying natural clusters, outliers, or technical batch effects without biasing the results toward pre-defined groups (see also Figure 1.8).

- o **Clustering:** Unsupervised clustering groups samples into subsets based on similarity, revealing hidden structure without outcome labels. Hierarchical clustering builds a tree (dendrogram) by iteratively merging the most similar pairs. K-means partitions samples into K groups by assigning each to its nearest centroid. Spectral clustering analyzes the eigen structure of a similarity graph to discover clusters of arbitrary shape. Chapter 12 covers these algorithms and their application to multi-omics integration in depth.
- o **Classification:** Supervised classification uses labeled training data to build models that predict group membership for new samples. Methods range from partial least squares discriminant analysis (PLS-DA) to more flexible approaches like random forests and support vector machines (SVM). The key to trustworthy classification is rigorous validation through cross-validation or independent test sets, guarding against overfitting (see Section 1.7).
- **Visualization ("Seeing is Believing"):** High-dimensional data requires visual synthesis to be interpretable. To transform abstract data matrices into intuitive biological maps, omics data science generally relies on three foundational categories of visualization.

- o **Scatter plots**: These project complex data points onto two-dimensional axes to reveal broad relationships. The two most critical scatter plots in omics are:
  - *PCA plots*: Essential for revealing global sample groupings, assessing batch effects, and identifying hidden outliers.
  - *Volcano plots*: Provide a global view of differential expression by plotting statistical significance (p-value) against biological effect size (fold change) for all features simultaneously.
- o **Heatmaps**: By representing numerical values as color gradients, heatmaps allow for the direct, simultaneous visualization of expression patterns across both samples and features. They are almost always paired with hierarchical clustering to group similar profiles.
- o **Networks**: Moving beyond individual feature abundances, these graphs map the functional relationships, biological pathways, and physical interactions (edges) between significant features (nodes) to reveal system-wide mechanisms.

**[DEFINITION: Heatmaps and Clustering]** A heatmap is a color-coded matrix that allows for the simultaneous visualization of complex omics data. Rather than plotting an entire dataset—which would be visually overwhelming—heatmaps typically display a targeted subset of highly variable or statistically significant features. Rows generally represent these features (genes, proteins, or metabolites), while columns represent individual samples. Each cell's color reflects the measurement's relative abundance, often utilizing a divergent color scale (e.g., red for up-regulation, blue for down-regulation) to make biological shifts immediately apparent. To enhance this view, clustering algorithms are applied to reorder the rows and columns based on their similarity. By grouping co-regulated features together horizontally and similar samples together vertically, the resulting clustered heatmap (Figure 1.5, panel 3) creates an information-dense visualization that allows you to see, at a single glance, which molecular signatures define specific sample groups and where potential outliers or sub-populations exist.

*Figure 1.5 — Key visualization types in omics. Four panels: (1) PCA score plot showing samples colored by condition, with clear separation between groups. (2) Volcano plot with $-\log_{10}$(p-value) on the y-axis and $\log_2$(fold-change) on the x-axis, highlighting significant features. (3) Heatmap of top significant features across samples, with hierarchical clustering. (4) Network diagram showing relationships between significant features.*

## Step 4: Functional Interpretation (Biology-Specific)

The final step connects statistical findings to biological mechanisms and generates testable hypotheses. A list of hundreds of differentially expressed genes is not, by itself, a biological insight. Functional interpretation transforms that list into a biological narrative by mapping features to known pathways, testing for enrichment in functional categories, and visualizing relationships in biological networks.

Chapter 2 covers functional interpretation in depth. While the statistical methods are the same or very similar across omics types, the function libraries often differ by layer (gene set libraries, metabolite set libraries, taxon set libraries, etc.).

The four-step workflow described in this section — raw data preprocessing, normalization, statistical analysis, and functional interpretation — provides the organizing structure for every discipline-specific chapter that follows.

*Figure 1.6 — Four-step analytical workflow. A horizontal flow diagram showing: Step 1 (Raw Data Preprocessing, technology-specific) → Step 2 (Normalization, omics-specific) → Step 3 (Statistical Analysis, omics-generic) → Step 4 (Functional Interpretation, biology-specific). Below each step: input format, key methods, and output. An arrow labeled "Chapter 2" points to Steps 1 and 4; the label "This chapter" points to Steps 2 and 3.*

# 1.6 Statistics, Machine Learning, and Artificial Intelligence

In the omics field, three analytical paradigms overlap and complement each other. Understanding when to use each is one of the most practically important skills in omics data science. The relationship between statistics, machine learning (ML), and artificial intelligence (AI) can be understood through their sample size requirements and analytical goals:

**Table 1.1 — Approach Overview**

| Approach | Sample Size | Best For | Examples | Trade-offs |
|---|---|---|---|---|
| **Statistics** | 10s–100s | Hypothesis testing, interpretability | t-tests, ANOVA, correlation, linear regression | Simple models, clear $p$-values, limited feature handling |
| **Machine Learning** | 100s–1,000s | Pattern recognition, prediction | Random forests, SVM, neural networks | Requires more data, harder to interpret, risk of overfitting |
| **Artificial Intelligence** | 1,000s+ | Complex tasks, automation, reasoning | Large language models, deep learning, generative models | Massive data and compute needs, interpretability challenges |

For the vast majority of omics studies — where sample sizes range from 10 to a few hundred — statistics and basic machine learning are the appropriate tools. Deep learning and AI approaches require thousands of samples to train reliably. This is why they are currently more useful for analyzing very large public datasets than for a typical laboratory study.

> **[TIP: Key Takeaway]** Statistics, ML, and AI differ primarily in their data requirements and interpretability. For typical omics studies with tens to hundreds of samples, classical statistics paired with careful visualization remains the most reliable and interpretable approach. Do not apply machine learning or deep learning to a study with less than 20 samples per group. The model will memorize your data rather than learn generalizable patterns. Start with simple statistical tests, and escalate to machine learning only when your sample size and analytical question justify it.

## A Critical Distinction: Significance vs. Performance

The distinction between **statistical significance** and **model performance** is a critical concept for omics researchers. These are fundamentally different questions that require distinct analytical approaches. Conflating them is a frequent source of error in interpreting high-dimensional data.

- **Significance: Is the effect real?** Measured by $p$-values and confidence intervals, significance answers: *"Is there evidence that this feature is associated with the outcome?"* It tells you if an effect is unlikely to have occurred by chance, but it does not tell you if that effect is large enough to be biologically or clinically useful.

- **Performance: Is the effect useful?** Measured by accuracy, sensitivity, specificity, and the Area Under the ROC Curve (AUC), performance answers: *"How well does this model predict the outcome in new samples?"* High performance indicates that the molecular signature reliably classifies unseen data.

A feature can be highly significant ($p < 0.001$) yet possess negligible predictive power because the actual change is too small to separate groups. Conversely, a combination of features—none of which are individually significant—can collectively drive excellent predictions. This is why robust omics analysis typically pairs both: we use significance to identify candidates and performance metrics to evaluate their collective utility.

The following contrast makes this concrete: a single feature with high statistical significance but negligible fold change contributes little to prediction, while a panel of modest features can collectively achieve strong classification accuracy.

*Figure 1.7 — Significance vs. Performance. Left: Volcano plot highlighting a feature with p = 0.001 but fold change of 1.05 (statistically significant but biologically trivial). Right: ROC curve showing that a panel of 10 features, none individually significant at p < 0.05, achieves AUC = 0.92 when combined. The message: significance and performance are different questions.*

# 1.7 Understanding Batch Effects and Technical Variation

**Batch effects** are among the most common and insidious sources of spurious results in omics studies. They arise whenever samples are processed in separate groups — different days, different reagent lots, different instruments, different operators — and the processing conditions introduce systematic differences that have nothing to do with biology. The danger is that batch effects can be *larger* than biological effects, leading researchers to discover "processing date" rather than "disease".

> **[DEFINITION: Batch Effect]** Systematic, non-biological variation arising from technical differences in sample processing, often from different processing dates, instruments, reagent lots, or operators.

## Identifying Batch Effects

Visual inspection is the most reliable way to detect batch effects. Three visualization approaches are particularly informative:

- **Distribution plots:** Compare the distribution of feature abundances across batches using box plots or density curves. Batches should have similar overall distributions if technical variation is minimal.

- **PCA plots:** Color samples by processing batch (date, instrument, operator) and check whether they separate by batch rather than by biological condition. If your PCA plot shows samples clustering by processing date, you have a batch effect.

- **Hierarchical clustering:** Similar to PCA, you can use a sample-wise dendrogram to check for technical grouping. If the primary branches split the samples perfectly by the day they were run rather than by their disease status, it provides clear secondary confirmation of a batch artifact.

## Addressing Batch Effects

Two fundamentally different strategies exist for handling batch effects:

**1. Batch effect removal (prior to analysis).** Computational methods like the *ComBat* algorithm estimate and remove the systematic variation attributable to batch identity, producing a "corrected" feature table. Internal standard normalization is another approach, where reference compounds or genes measured alongside your samples are used to calibrate across batches. The corrected table is then analyzed normally.

**2. Batch modeling (during analysis).** Rather than removing batch effects from the data, some methods allow you to include batch as a **covariate** in your statistical model. For example, the model "outcome ~ treatment + batch" estimates the treatment effect while accounting for batch differences. This approach is more statistically principled because it does not alter the raw data.

> **[CAUTION: Don't Ignore Batch Effects]** Unaddressed batch effects generate false discoveries. Always inspect your data for batch effects before running statistical analyses. If your treatment groups are perfectly confounded with batches (all treatments processed on Day 1, all controls on Day 2), no statistical method can separate the biological effect from the batch effect — the experiment must be redesigned.

The following panels contrast a problematic batch-confounded PCA with a corrected version, and illustrate how balanced experimental design prevents confounding in the first place.

*Figure 1.8 — Batch effect detection and correction. Top-left: PCA plot where samples separate by processing batch rather than biological condition (problematic). Top-right: same data after ComBat correction, now separating by condition. Bottom: experimental design showing a confounded design (treatment = batch) versus a balanced design (treatment groups split across batches).*

## 1.8 Classification, Overfitting, and Cross-Validation

When the goal is to build a model that classifies samples — predicting disease versus healthy, responder versus non-responder — **overfitting** becomes the central concern. Understanding why overfitting occurs in omics data, and how to prevent it, is essential for any supervised analysis.

> **[DEFINITION: Overfitting]** When a model learns the training data too well, capturing noise rather than true biological patterns. Results appear excellent on training data but fail to generalize to new, unseen samples.

## Why Omics Data Is Prone to Overfitting

The $p \gg n$ problem makes overfitting almost inevitable when building classifiers on omics data without proper safeguards. With thousands of features and few samples, the model has enormous flexibility to find patterns — including patterns that are purely due to noise. By chance alone, some features will correlate with the outcome in your specific dataset. A classifier that uses these chance correlations will perform brilliantly on the training data but fail completely on new samples.

This is one of the most common failures in published omics studies. Researchers build a classifier, report excellent accuracy on their training data, and the model fails to replicate in independent cohorts. The root cause is almost always overfitting.

## The Solution: Cross-Validation

**Cross-validation** provides unbiased estimates of how a model will perform on new data. The principle is straightforward: partition your data into separate training and testing portions, build the model on the training portion, and evaluate it on the testing portion. Since the test data was not used to build the model, performance on the test data reflects how the model will generalize.

The most common approach is **k-fold cross-validation**: the data is split into $k$ equal portions (folds). For each fold, the model is trained on the remaining $k-1$ folds and tested on the held-out fold. After all folds have served as the test set, the average performance across folds provides a realistic estimate of generalization accuracy.

> **[BEST PRACTICE: Always Use Cross-Validation]** Never report accuracy or performance metrics from your training data alone. Always use independent test data or cross-validation to get unbiased performance estimates. A model that achieves 95% accuracy on training data but 55% on cross-validation is overfitting — the 55% figure is the honest estimate of performance.

The following schematic shows how 5-fold cross-validation rotates through the data so that every sample serves exactly once as a test case, producing five independent performance estimates whose average reflects true generalization accuracy.

*Figure 1.9 — Cross-validation schematic. A diagram showing 5-fold cross-validation: data split into 5 portions, with each portion taking a turn as the test set (highlighted) while the remaining 4 are used for training. Five performance metrics are produced and averaged.*

# 1.9 The Feature Table: Foundation of Omics Analysis

The feature table is the universal starting point for omics statistical analysis. Regardless of whether your data comes from RNA-seq, mass spectrometry, or 16S amplicon sequencing, the processed data takes the form of a matrix.

## Structure and Conventions

By convention in this book, the feature table follows a consistent orientation:

- **Rows = Features:** Genes, proteins, metabolites/peaks, or microbial taxa. Each row is one measured entity.

- **Columns = Samples:** Biological replicates or individual subjects. Each column is one sample.

- **Cells = Abundance values:** Expression levels, peak areas, read counts, or relative abundances. These are the quantitative measurements.

Feature tables from different technologies look different on the surface. RNA-seq tables contain integer counts with gene IDs like ENSG00000141510, while untargeted metabolomics tables contain continuous intensities with peak IDs like 401.05_2.3. Despite their surface differences, the analytical principles that apply to them are the same. This is one of the key insights of this book: the same statistical and visualization methods apply across omics types, because they all share the same tabular structure.

*Figure 1.10 — Feature table examples. Two tables side by side: Left: RNA-seq count table with gene IDs (ENSG001, ENSG002, …) as rows and sample IDs as columns, with integer count values. Right: LC-MS feature table with m/z_RT IDs as rows and sample IDs as columns, with continuous intensity values. A bracket labeled "Same analytical framework" spans both tables.*

## Metadata: The Essential Context

A feature table without metadata is a matrix of meaningless numbers. **Metadata** describes the biological and technical context of each sample: clinical diagnosis, treatment assignment, demographics (age, sex), processing information (batch, date, operator), and quality metrics. Without metadata, you cannot perform any hypothesis-driven analysis.

Metadata can be organized in two ways:

- **Simple format:** Metadata columns are included alongside the feature table in a single file. This works well for small studies with one or two metadata variables.

- **Complex format:** The feature table and metadata are stored as separate files, linked by sample identifiers. This is preferred for studies with complex experimental designs and many metadata variables.

> **[BEST PRACTICE: Metadata Best Practices]** Use clear, consistent naming without special characters. Always include batch and processing information, even if you do not plan to analyze it. Document your data collection procedures. Many analyses fail not because of bad omics data but because of incomplete metadata.

## Two Fundamental Operations on Feature Tables

Every analysis of a feature table approaches the data from one of two complementary perspectives. To build a complete biological story, a researcher must move fluently between these two "spaces".

- **The Feature-Centric Perspective (Univariate Analysis):** This perspective examines one feature at a time across all samples. For each gene, protein, or metabolite, you test for an association with the outcome to produce a $p$-value and an effect size. While highly interpretable—you know exactly which features are significant—this approach ignores interactions between variables and requires rigorous multiple testing correction to avoid false discoveries.
- **The Sample-Centric Perspective (Multivariate Analysis):** This perspective considers all features simultaneously, treating each sample as a single point in a high-dimensional space. Methods like PCA and clustering reveal how samples relate to one another and what global patterns exist. While this captures the complex, interconnected structure of biology, it can be mathematically challenging to trace these global shifts back to specific, individual driving features.

Neither perspective is sufficient on its own. Univariate analysis identifies the individual "actors" (significant features), while multivariate analysis reveals the "plot" (the overall structure). Together, they validate that your individual findings are part of a coherent, system-wide pattern rather than isolated, statistical noise.

## 1.10 The XiaLab Tool Ecosystem

The XiaLab at McGill University has developed an integrated suite of web-based tools for omics data analysis. These tools implement the four-step analytical framework described in this

chapter and provide accessible interfaces that require no programming — while maintaining the statistical rigor needed for publication-quality results.

**Table 1.2 — Core Tools**

| Tool | Primary Use | URL Address |
|------|-------------|-------------|
| **ExpressAnalyst** | Gene expression analysis | www.expressanalyst.ca |
| **ProteoAnalyst** | Quantitative proteomics analysis | www.proteoanalyst.ca |
| **MetaboAnalyst** | Metabolomics analysis | www.metaboanalyst.ca |
| **MicrobiomeAnalyst** | Microbiome community analysis | www.microbiomeanalyst.ca |
| **NetworkAnalyst** | Network analysis and visualization | www.networkanalyst.ca |
| **OmicsNet** | Multi-omics network integration | www.omicsnet.ca |
| **OmicsAnalyst** | Multi-omics data integration | www.omicsanalyst.ca |
| **miRNet** | MicroRNA network analysis | www.mirnet.ca |

All XiaLab tools share a common design philosophy: sensible defaults that produce reliable results for most analyses, with expert options available for specialized needs. Data upload accepts simple file formats (CSV, TSV), visualization is interactive with publication-quality export, and comprehensive tutorials and video walkthroughs are available for each tool.

In addition to supporting standardized raw data formats for NGS (FASTQ) and MS (mzML), as well as outputs from well-established processing pipelines, XiaLab tools accept common tabular formats. Two organizational approaches are supported:

- **Simple format:** Metadata columns (sample ID, group label, covariates) are included alongside the feature abundance data in a single file. This is convenient for small studies with one or two metadata variables.
- **Complex format:** The feature table and metadata are stored as separate files, linked by sample identifiers. This is preferred for studies with complex experimental designs and many metadata variables.

> **[TIP: Invest Time in Tutorials]** For subscribed users, each XiaLab tool has excellent slides and video tutorials. Invest 1–2 hours learning the interface

> before analyzing your own data. This prevents common mistakes and accelerates your analysis significantly.

The following diagram maps the XiaLab tools to their respective data types and shows how results from one tool can flow into another.

*Figure 1.11 — The XiaLab tool ecosystem. A diagram showing all tools arranged by data type (columns: gene expression, metabolomics, microbiome, multi-omics) with arrows showing how results can flow between tools.*

•   •   •   •   •

## Pro Omics Tool Suite: Enhanced Features for Professional Research

To sustain our community tools and better support research groups requiring persistent project management, automated reporting, and extended analytical modules, XiaLab provides the **Pro Omics Tool Suite** — a professional-tier platform that spans six integrated tools: *ExpressAnalyst*, *MetaboAnalyst*, *MicrobiomeAnalyst*, *miRNet*, *OmicsNet*, and *OmicsAnalyst*. All six Pro tools share the following capabilities.

## Addressing Reproducibility Through Project Management

One of the most practical challenges is maintaining continuity across analysis sessions. A typical omics study involves dozens of decisions — filtering thresholds, normalization options, statistical tool, significance cutoffs, enrichment databases — and revisiting an analysis months later during manuscript revision often means repeating the entire workflow from scratch.

The Pro Omics Tool Suite addresses this with **project management**. Every analysis session can be saved, organized into projects, and resumed at any point. The platform maintains a complete record of all parameters, intermediate results, and figures generated during an analysis. When a reviewer asks you to re-run the analysis with a different normalization method or a stricter FDR threshold, you can reload your project, make the specific change, and regenerate results without repeating upstream steps. This persistence also supports collaborative work: team members can access shared project storage to review or extend each other's analyses.

## Ensuring Methodological Consistency Through Automated Workflows

Beyond saving individual sessions, the Pro tools include a **workflow automation system** that captures multi-step analysis pipelines as reusable templates. A workflow is represented as a directed acyclic graph (DAG) of analysis steps — from data upload through normalization, statistical analysis, and functional interpretation — with each step's parameters recorded. Once

defined, a workflow can be applied to new datasets, ensuring that the same analytical protocol is applied consistently across studies.

This feature is particularly valuable for research group that wants to ensure methodological consistency for studies involving omics data, core facilities processing multiple datasets with standardized protocols; and multi-site studies requiring identical analytical pipelines.

## Accelerating Publication Through Automated Reporting

Each Pro tool can generate comprehensive analysis reports automatically. With a single click, the platform compiles all figures, statistical tables, methods descriptions, and parameter documentation into a formatted HTML or PDF report. The report includes every visualization generated during the analysis — from QC plots through statistical summaries and enrichment networks — along with the statistical details needed for a methods section.

For presentations, the Pro tools offer two output formats: *Reveal.js* — a **web-based presentation** framework that produces interactive, browser-ready slide decks ideal for online sharing and remote presentations — and **PowerPoint** for traditional offline use. Both formats extract key figures and results into presentation-ready slides, eliminating the manual process of exporting individual figures, adjusting formatting, and assembling them by hand. When an analysis is revised, the report and slide deck can be regenerated instantly to reflect the latest results.

## Breaking Interface Constraints with AI-Powered Assistance

The Pro Omics Tool Suite integrates an AI assistant that provides context-aware guidance throughout the analysis. Each tool's assistant is specialized in its respective domain, drawing on a curated knowledge base of common questions, best practices, and troubleshooting advice. Users can ask questions in natural language—such as "Which normalization method should I use for my RNA-seq data with batch effects?" or "How do I interpret this volcano plot?"—and receive targeted guidance directly within the platform.

Beyond general support, the AI powers natural-language plot customization. Rather than navigating nested parameter menus to adjust visualizations, you can describe desired changes in plain language (e.g., "Increase the point size and highlight significant genes in red"). The system instantly translates these instructions into the appropriate plotting parameters, lowering the technical barrier for researchers who need publication-quality figures but are not comfortable with R syntax.

Most importantly, this AI integration expands the platform's analytical boundaries by transcending the inherent constraints of a graphical user interface (GUI). Web-based tools traditionally restrict researchers to pre-programmed menus. However, if your study requires a bespoke statistical model or a specialized data transformation not natively built into the interface, you can request it conversationally. The AI translates these complex analytical concepts into custom code executed seamlessly in the background (*Note: Because executing AI-generated code requires stringent security and reproducibility controls, this feature is currently undergoing rigorous validation and sandboxing prior to its official release*).

## Ensuring Data Sovereignty Through Flexible Deployment

While the public versions of the XiaLab tools run exclusively on cloud servers, the Pro Omics Tool Suite offers flexible deployment. Researchers can choose cloud-based analysis for convenience or install the platform locally for full control over data and computing resources. Local installation is essential for studies involving sensitive human genomic data subject to institutional review board (IRB) restrictions or data governance policies (such as HIPAA or GDPR) that prohibit uploading to external servers. It also enables organizations with dedicated computing infrastructure to leverage their own hardware for large-scale analyses.

> **[TIP: Choosing Between Public and Pro Tiers]** While the public XiaLab tools are excellent for learning, initial data exploration, and analyzing small-to-medium datasets, the Pro Omics Tool Suite is highly recommended for routine, rigorous research. Because the Pro tier runs on dedicated, distributed cloud instances (or can be securely installed on-premise), it guarantees the computational power and stability required for high-throughput analyses. Importantly, subscribing to the Pro tier directly helps sustain the free community versions, ensuring these foundational tools remain accessible to researchers worldwide

## Summary

This chapter has established the conceptual framework that underlies all omics data analysis. The key principles are:

- Omics technologies measure biological molecules at multiple layers — genome, transcriptome, proteome, metabolome — each providing complementary information about biological systems.

- Two distinct challenges define omics data analysis: the **size challenge** (managing massive raw data files, solved through standardized pipelines) and the **complexity challenge** (analyzing high-dimensional feature tables, requiring scientific judgment and multiple analytical approaches).

- The **data distillation workflow** applies across all omics types. It progressively reduces ~100,000 raw features to fewer than 5 actionable biological insights through data processing, statistical testing, and functional interpretation.

- Statistics, machine learning, and AI form a continuum of analytical approaches, each suited to different sample sizes and analytical goals. For typical omics studies, statistics and basic machine learning are the appropriate tools.

- Batch effects, overfitting, and multiple testing are the three most common pitfalls. Visual inspection, cross-validation, and FDR correction are the essential safeguards.

- The **feature table** — a matrix of features × samples × abundance values — is the universal starting point for omics analysis, regardless of the underlying technology.

- The **Pro Omics Tool Suite** extends the public tools with advanced features designed for professional research workflows: persistent project management, automated workflow engines, presentation-ready analysis reports and slide decks, an integrated AI assistant for custom coding, and flexible on-premise deployment for data governance compliance.

With this conceptual foundation in place, Chapter 2 dives into the practical details at the two ends of the workflow: how raw data from sequencing instruments and mass spectrometers is transformed into feature tables (Step 1), and how lists of significant features are interpreted through enrichment analysis and network context (Step 4).

•   •   •   •   •

**Table 1.3 — Key Terms**

| Term | Definition |
| --- | --- |
| **Omics** | The comprehensive, high-throughput study of a complete set of biological molecules within a system. |
| **Feature Table** | Data matrix with features (rows) × samples (columns) × abundance values; the universal input for omics analysis. |
| **Data Distillation** | The process of progressively reducing high-dimensional data to actionable biological insights through data processing, filtering, statistics, and knowledge integration. |
| **Normalization** | Statistical adjustment for technical variation (sequencing depth, sample amount, instrument sensitivity), ensuring across-sample comparability and approximating a normal distribution. |
| **Batch Effect** | Systematic non-biological variation from technical differences in sample processing (date, instrument, reagent lot, operator). |
| *p*-**value** | The probability of observing a result as extreme as the data under the null hypothesis; describes random chance, not biological strength. |
| **Effect Size** | The magnitude of biological change (e.g., $\log_2$ fold-change, Cohen's $d$); complements *p*-values. |
| **FDR (False Discovery Rate)** | The expected proportion of false positives among all features declared significant; controlled by Benjamini–Hochberg correction. |
| **Multiple Testing Correction** | Statistical adjustment for the increased false-positive risk when testing thousands of hypotheses simultaneously. |
| **Overfitting** | When a model captures noise rather than true patterns, performing well on training data but failing on new samples. |
| **Cross-Validation** | A method for estimating model performance by repeatedly training on subsets and testing on held-out data. |
| **PCA (Principal Component Analysis)** | An unsupervised method that identifies axes of maximum variance in high-dimensional data for visualization and dimension reduction. |
| **Supervised Analysis** | Methods that use outcome labels to select features or build predictive models. |
| **Unsupervised Analysis** | Methods that identify patterns without using outcome labels (PCA, clustering); no risk of overfitting to outcomes. |

| Term | Definition |
|------|-----------|
| **Semi-Supervised Analysis** | Combining data-driven results with biological knowledge for robust interpretation. |

• • • • •

## Review Questions

**1.** Explain the difference between the size challenge and the complexity challenge in omics data analysis. Why do they require different solutions?

**2.** A researcher normalizes their RNA-seq data using RPKM and then provides these values to *DESeq2* for differential expression analysis. Why is this likely to produce incorrect results?

**3.** You perform PCA on a metabolomics dataset and find that samples separate primarily by processing date rather than disease status. What does this indicate, and what strategies could you use to address it?

**4.** A colleague reports that their machine learning classifier achieves 98% accuracy on a dataset with 30 samples per group and 15,000 features. Why should you be skeptical of this result, and what validation approach would you recommend?

**5.** Explain why multiple testing correction is necessary in omics studies. If you test 10,000 features at $p < 0.05$ without correction, how many false positives would you expect by chance alone?

**6.** Why is it important to combine both univariate (feature space) and multivariate (sample space) analyses when exploring a feature table? What does each approach reveal that the other misses?

**7.** A metabolomics study identifies 500 significant peaks. Explain why the data distillation process does not stop here, and what additional steps are needed to arrive at biological insights.

• • • • •

## Further Reading

- Ritchie, M.E., *et al.* "Methods of integrating data to uncover genotype–phenotype interactions." *Nature Reviews Genetics* 16 (2015): 85–97.
- Hasin, Y., Seldin, M., and Lusis, A. "Multi-omics approaches to disease." *Genome Biology* 18, 83 (2017).

- Xia, J., and Wishart, D.S. "Using *MetaboAnalyst* 3.0 for comprehensive metabolomics data analysis." *Current Protocols in Bioinformatics* 55 (2016): 14.10.1–14.10.91.

- Pang, Z., *et al.* "*MetaboAnalyst* 6.0: narrowing the gap between raw spectra and functional insights." *Nucleic Acids Research* 52, W1 (2024).

- Leek, J.T., *et al.* "Tackling the widespread and critical impact of batch effects in high-throughput data." *Nature Reviews Genetics* 11 (2010): 733–739.

- Noble, W.S. "How does multiple testing correction work?" *Nature Biotechnology* 27, 12 (2009): 1135–1137.

- Hastie, T., Tibshirani, R., and Friedman, J. *The Elements of Statistical Learning.* 2nd ed. Springer, 2009.

- Zhou, G., *et al.* "*OmicsAnalyst*: a comprehensive web-based platform for visual analytics of multi-omics data." *Nucleic Acids Research* 49, W1 (2021): W476–W482.

## What's Next

Chapter 2 fills in the practical details for the first and last steps of the four-step workflow. It begins with the technology-driven data paradigms (NGS-omics, MS-omics, spatial omics) and walks through how raw sequencing data and mass spectrometry data are processed into structured feature tables. It then covers functional interpretation — enrichment analysis, pathway topology, and network context — showing how lists of significant features are transformed into biological narratives.

# Chapter 2: From Raw Data to Biological Functions

## Learning Objectives

By the end of this chapter, you will be able to:

- Distinguish between the two technology-driven data paradigms in omics (NGS-omics, MS-omics) and explain how each transforms raw measurements into feature tables.

- Describe the NGS raw data processing pipeline — from FASTQ files through quality control, read mapping, and assembly to raw count tables — including the role of k-mers.

- Explain the LC-MS workflow — from chromatographic separation through ionization and mass detection to peak tables — and articulate the difference between relative and absolute quantification.

- Compare and contrast Over-Representation Analysis (ORA) and Gene Set Enrichment Analysis (*GSEA*) as strategies for enrichment analysis, including their assumptions and potential biases.

- Explain how pathway and network context enhance functional interpretation beyond simple enrichment tests.

• • • • •

## Introduction

Chapter 1 introduced two fundamental challenges in omics data analysis: the size challenge posed by raw instrument data, and the complexity challenge that emerges once we have a structured feature table. It also outlined a four-step analytical workflow — from raw data preprocessing through data processing, statistical analysis, and functional interpretation — and demonstrated that this workflow applies across all omics disciplines.

This chapter fills in the two ends of that workflow. We begin with **Step 1**: how raw instrument output is transformed into a structured feature table. This step is inherently technology-specific — sequencing data and mass spectrometry data require entirely different computational pipelines. We then move to **Step 4**: how a list of significant features is interpreted through the lens of biological function using enrichment analysis, pathway analysis, and network context.

The logic follows the data distillation funnel introduced in Chapter 1. At the top of the funnel, researchers face gigabytes or even terabytes of vendor-specific raw files. By the time raw data

processing is complete, those files have been condensed into a feature table — a matrix of features by samples by abundance values — that is orders of magnitude smaller (typically kilobytes to megabytes) and ready for statistical analysis.

At the bottom of the funnel, after statistical and machine learning methods have identified significant features, functional analysis condenses thousands of individual molecules into a handful of enriched pathways. This gives researchers the biological narrative they need to generate hypotheses and design follow-up experiments.

Between these two ends lie the data processing, statistical analysis, and visualization steps introduced in Chapter 1. Those steps are largely omics-generic — they work the same way whether your features are genes, metabolites, or proteins. Raw data processing and functional interpretation, by contrast, require **technology-specific** and **biology-specific knowledge**. That is what this chapter provides.

## 2.1 Two Fundamental Technology-Driven Data Paradigms

Omics is fundamentally technology driven. Before diving into the mechanics of specific processing pipelines, it is crucial to recognize that almost all standard omics technologies fall into one of two fundamental data generation paradigms. Rather than being defined by the biological molecules under study, these categories are defined by the fundamental nature of the measurements the instrument produces. Understanding this binary clarifies why different omics disciplines require completely distinct computational approaches to ultimately converge on a unified feature table.

### 1. Sequence-Based Omics (NGS-omics)

NGS-omics uses next-generation sequencing to quantify biological molecules by counting high-throughput sequence reads. In this paradigm, discrete, digital counts of DNA or RNA fragments serve as the primary data unit. The core computational objective is to align these millions of short reads to a reference genome (or assemble them *de novo*) to construct a structured feature table.

Because reads are discrete, NGS data is inherently digital: a gene's expression level is typically proportional to the number of reads mapped to its locus. While this counting nature simplifies quantification by removing the need to model complex peak shapes, it introduces its own unique computational challenges—namely, accounting for vastly different sequencing depths across samples, relying on high-quality reference genomes, and managing the immense computational cost of processing billions of short text strings.

## 2. Mass Spectrometry-Based Omics (MS-omics)

MS-omics relies on measuring the **mass-to-charge ratio (m/z)** of ionized molecules. Unlike the digital counts of NGS, the raw data consists of continuous signal intensities and complex spectral peaks. The objective is to distill these "analog" signals into biological features — a task that differs by molecular layer:

- **In global untargeted metabolomics:** The focus is on **peak deconvolution**—the process of resolving overlapping signals, grouping related isotopes and adducts, and aligning features across samples to identify discrete small molecules.

- **In bottom-up quantitative proteomics:** The workflow requires identifying peptide sequences from fragment spectra (MS2), followed by the **statistical inference** of proteins from those constituent peptide sets.

A defining challenge in MS-omics is the many-to-many relationship between signals and molecules. A single analyte (whether a metabolite or a peptide) typically produces multiple peaks through isotopes, various charge states, and in-source fragmentation. Conversely, a single m/z window may harbor multiple co-eluting species that must be mathematically separated. Managing this resulting technical noise and the high proportion of missing values—often caused by signals falling near the instrument's limit of detection or the stochastic nature of ion selection—is the central goal of the MS data processing pipeline.

### Emerging and Integrative Paradigms

While sequence-based and mass-based measurements form the foundation of modern data science, two rapidly evolving paradigms build upon these core technologies to capture even greater biological complexity:

- **Spatial Omics:** This extension captures the molecular state of a sample while preserving its intricate physical architecture. Unlike traditional omics, which grinds tissue into a "bulk" average, spatial methods maintain the 2D or 3D physical coordinates of every signal. Depending on the platform, the raw data unit may be a pixel/voxel (in MS imaging) or a spatially barcoded spot (in NGS). The primary computational objective is to integrate high-dimensional molecular signals with morphological features, often requiring image segmentation to assign transcripts or proteins to specific physical coordinates.

- **Multi-Modal Omics:** This integrative approach seeks to fuse discrete sequence counts, continuous spectral intensities, and spatial coordinates into a single, unified analytical framework. The objective is to overcome the inherent blind spots of any single technology. For example, a researcher can explicitly correlate the presence of a specific

microbial sequence (NGS-omics) with the production of a localized metabolite (MS-omics) within a specific functional tissue niche (spatial omics). This cross-paradigm distillation provides a holistic view essential for understanding complex systems biology.

**Table 2.1 — Comparing the Paradigms**

| Omics Type | Data Unit | Processing Objective | Output |
|---|---|---|---|
| **NGS-omics** | Digital counts (reads) | Read alignment or *de novo* assembly | Gene/transcript count table |
| **MS-omics** | Continuous intensities (spectra) | Peak picking, alignment, identification | Quantified molecule table |
| **Spatial Omics** | Pixels/voxels, barcoded spots | Image segmentation, cell delineation | Spatially indexed table |
| **Multi-modal** | Mixed (counts + intensities + spatial) | Cross-paradigm fusion, latent space alignment | Unified multidimensional feature table |

The remainder of this chapter focuses strictly on the two dominant paradigms—NGS-omics and MS-omics—which together account for the vast majority of omics data generated today. Despite their fundamentally different measurement principles (digital counting versus analog signal processing), both paradigms ultimately converge on the exact same destination: a structured, quantitative **feature table**. Once the raw reads are aligned or the raw peaks are integrated, the underlying hardware ceases to matter as much. The resulting matrix is now ready for the statistical and functional analyses described throughout the rest of this book.

*Figure 2.1 — Schematic showing the three data paradigms side by side. Left: NGS-omics (reads mapping to genome, producing count table). Center: MS-omics (spectra showing peaks, producing intensity table). Right: Spatial omics (tissue image with segmented cells, producing spatial molecular map). All three converge on a feature table at the bottom.*

## 2.2 NGS Raw Data Processing

Next-generation sequencing has transformed biology by making it possible to read the nucleotide sequences of millions of DNA or RNA fragments in a single experiment. For most omics researchers, the goal is straightforward: start with raw sequencing files, apply a computational pipeline, and produce a table of gene counts that can be analyzed statistically. This section walks through each stage of that pipeline.

## The FASTQ File: Where It All Begins

When a sequencing center returns your data, you receive FASTQ files — often several gigabytes each. A **FASTQ** file stores both the nucleotide sequences (the reads) and their associated quality scores. Each read occupies exactly four lines:

- **Line 1:** A header beginning with @ that identifies the sequencing instrument, flow cell, and read coordinates.
- **Line 2:** The nucleotide sequence itself (e.g., ATCGATCG…).
- **Line 3:** A separator line (usually just a + sign).
- **Line 4:** Quality scores encoded as ASCII characters, one per base.

The four-line structure of a FASTQ read — header, sequence, separator, quality — is compact but information-rich, encoding both the nucleotide identity and the sequencer's confidence in each base call.

> [DEFINITION: FASTQ] The standard file format for storing nucleotide sequences together with their per-base quality scores from next-generation sequencing instruments. Each read occupies four lines: header, sequence, separator, and quality string. FASTQ files can be many gigabytes in size. Opening them in a text editor or spreadsheet program will freeze your computer. Use command-line tools like *head -n 16 myfile.fastq* to inspect just the first few reads.

*Figure 2.2 — Annotated FASTQ file showing four reads. Each read has four lines: header (@), sequence, separator (+), and quality scores. Annotations point to the instrument ID, flow cell ID, lane number, tile, and x/y coordinates in the header line.*

## Quality Scores and the Phred Scale

Each base in a read is assigned a quality score by the sequencing instrument, reflecting the machine's confidence in that particular base call. These scores use the **Phred scale**:

$$Q = -10 \times \log_{10}(P)$$

where *P* is the probability that the base call is incorrect. A Phred score of Q20 means a 1-in-100 chance of error (99% accuracy), while Q30 means 1 in 1,000 (99.9% accuracy). Modern Illumina platforms routinely produce data where the majority of bases exceed Q30.

**Table 2.2 — Phred Score Overview**

| Phred Score | Error Probability | Accuracy | Interpretation |
| --- | --- | --- | --- |
| Q10 | 1 in 10 | 90% | Poor quality |
| Q20 | 1 in 100 | 99% | Acceptable minimum |
| Q30 | 1 in 1,000 | 99.9% | High quality (modern standard) |
| Q40 | 1 in 10,000 | 99.99% | Excellent |

Quality scores are encoded as ASCII characters rather than numbers to save space. Each character maps to a quality value, ensuring that each position in the quality string aligns perfectly with the corresponding base in the sequence. You do not need to decode these characters manually — quality control tools handle this automatically.

## Quality Control: *FastQC* and *fastp*

Before any downstream analysis, you should assess the quality of your sequencing data. Tools like *FastQC* provide a visual summary of quality metrics across all reads, including per-base quality score distributions, GC content, adapter contamination, and sequence duplication levels.

In the early days of NGS, quality scores ($Q$) often dropped sharply toward the 3' end of reads. This necessitated aggressive trimming, or the removal of low-quality bases beyond a specific threshold to prevent false-positive alignments. While modern Illumina data is usually of such high quality that aggressive trimming is rarely required, a thorough quality check remains essential to catch issues like adapter contamination or failed sequencing lanes. Currently, *fastp* has become the industry standard for preprocessing; it integrates quality trimming, adapter removal, and quality reporting into a single, high-performance step. For most modern datasets, the default parameters in fastp are sufficient to prepare the data for alignment.

> **[TIP: Quality Is Rarely the Bottleneck]** Modern sequencing platforms produce consistently high-quality data. The real bottleneck in NGS data analysis is the computational cost of aligning billions of reads, not the quality of those reads. Quality control is a quick sanity check, not the rate-limiting step.

The contrast between early-era and modern sequencing quality is dramatic: older platforms showed steep quality degradation at read ends, while current Illumina instruments maintain high confidence across the full read length.

## The K-mer Concept

Before discussing read mapping and assembly, it is worth understanding a concept that underpins both: the **k-mer**. A k-mer is a contiguous subsequence of length $k$ extracted from a read using a sliding window. For example, the sequence ATCGATCG contains the following 4-mers: ATCG, TCGA, CGAT, GATC, ATCG.

Why are k-mers useful? A single nucleotide carries very little information — it could match millions of positions in a genome. But as $k$ increases, the k-mer becomes increasingly unique. In *E. coli*'s small genome, k-mers of length 10–12 are sufficient to map uniquely to a single location. In the much larger human genome, uniqueness requires approximately k = 25–32. This means that short reads of 75–150 bases contain more than enough information for unique mapping — we do not need extremely long reads for gene expression quantification.

> **[DEFINITION: K-mers as Dimensionality Reduction]** The shift from analyzing individual bases to analyzing k-mers is a form of dimension reduction. K-mers capture contextual sequence information in a computationally efficient representation. This concept parallels *n-grams* in natural language processing, where combinations of words (bigrams, trigrams) carry more meaning than individual words. K-mers are widely used in genome assembly, read mapping, and error correction. They serve as powerful features for machine learning, allowing algorithms to detect motifs and structural patterns without the computational burden of full sequence alignment.

K-mers serve three critical roles in NGS data processing:

- **Read mapping:** K-mer-based indexing allows alignment algorithms to quickly locate candidate positions in a reference genome, dramatically reducing search time.
- **Genome assembly:** Overlapping k-mers from different reads are used to reconstruct contiguous sequences (contigs). By adjusting $k$, assemblers can balance sensitivity (smaller $k$ captures more overlaps) against specificity (larger $k$ avoids false joins in repetitive regions).

- **Sequence signatures:** The frequency distribution of k-mers in a genome or metagenome serves as a compact signature of its composition, useful for taxonomic classification and sequence comparison.

The relationship between k-mer length and genome complexity determines how long a subsequence must be to achieve unique mapping — a principle that explains why short reads are fully adequate for expression quantification in well-annotated genomes.

*Figure 2.4 — K-mer uniqueness versus genome complexity. A chart showing that k-mers of length ~12 are sufficient for unique mapping in E. coli, while k = 25–32 is needed for the human genome. Illustrates why short reads (75–150 bp) contain enough information for gene expression mapping.*

# 2.3 Quantification: Transforming Reads to Gene Counts

The primary goal is to transform raw sequence reads into a structured gene count table. For organisms with well-characterized reference genomes—including humans, mice, and reference-grade bacteria—the standard approach is to map reads to a reference. This process identifies the origin of each read, though the specific computational strategy depends on the biological architecture of the target and the required analysis speed.

## *Traditional Coordinate-Based Alignment*

Traditional aligners perform a base-by-base search to "pin" a read to a specific location:

- **Bacterial/genomic mapping:** *Bowtie2* is the gold standard for microbial data and genomic DNA. Because bacteria lack introns, *Bowtie2* uses a high-speed, "unspliced" logic that is ideal for the linear gene structures of microbiota.
- **Eukaryotic transcriptomics:** *HISAT2* and *STAR* are preferred for host RNA-seq because they are **splice-aware**. They can handle reads that span exon–exon junctions, accurately mapping sequences where the intervening intron has been removed.

## *Direct Transcriptome Inference (Alignment-Free)*

While traditional aligners map reads to a genomic location and then "count" overlaps, **pseudo-aligners** like *Kallisto* and *Salmon* infer the transcriptome directly. By breaking reads into k-mers (short, fixed-length sequences) and matching them against a pre-built transcriptome index, these tools skip the computationally intensive step of calculating exact genomic coordinates. This "alignment-free" approach is highly efficient for standard RNA-seq, where the primary goal is the rapid quantification of known genes rather than the discovery of novel genomic architectures or unannotated splice sites.

The schematic below traces reads from their raw sequenced form through alignment to a reference genome, illustrating the handling of uniquely mapped reads, splice-junction reads, and multi-mapped reads, and showing how each contributes to the final per-gene count.

*Figure 2.5 — Read mapping schematic. Short reads aligned to a reference genome showing: (1) uniquely mapped reads spanning a single gene, (2) reads spanning an exon–exon junction (splice junction), (3) multi-mapped reads in a repetitive region, and (4) the resulting raw count per gene.*

## Strategies for Non-Model Species

Not every study involves an organism with a high-quality reference genome. For non-model species—such as many fish, amphibians, and plants—researchers cannot use standard coordinate-based mapping. Instead, they must employ strategies to bridge the gap between raw reads and a functional feature table.

- *De Novo* **Assembly (Reconstruction)**: This strategy involves reconstructing the genome or transcriptome directly from the sequencing reads without a template. Standard tools like *SPAdes* use k-mer-based algorithms to organize overlapping fragments into a *de Bruijn graph*, which is then traversed to build contiguous sequences called contigs. While this is the only way to discover the unique genomic architecture of a new species, it is computationally intensive and becomes exponentially more difficult for large, complex, or polyploid genomes.

- **Ortholog Mapping (Functional Shortcut)**: For researchers primarily interested in gene expression or pathway analysis rather than genome structure, ortholog mapping provides a faster alternative. Instead of assembling a new reference, tools like *Seq2Fun* map reads directly against a database of conserved protein-coding genes from closely related species. This effectively bypasses the assembly phase, allowing for rapid functional quantification of the "known" protein-coding space within a non-model organism.

> **[TIP: Know Your Goal]** If your goal is to quantify gene expression in a non-model species, you may not need a full genome assembly. Ortholog-based approaches like Seq2Fun can produce reliable expression tables directly from raw reads, saving months of assembly and annotation work.

## Generating the Raw Count Table

The final stage of the NGS primary pipeline is the conversion of processed sequence data into a quantitative matrix. Whether using traditional mapping or alignment-free inference, the goal is to "assign" each read to a specific biological feature.

- **In traditional coordinate-based workflows**, the mapping stage is followed by quantification using tools like *featureCounts* or *HTSeq*. These tools apply rigorous rules to "distill" mapped alignments into a count table, determining exactly which gene or feature should be credited for each read. This stage must handle complex biological scenarios, such as reads that map to antisense strands or regions where multiple genes overlap on the same genomic locus.

- **In alignment-free workflows**, tools like *Kallisto* or *Salmon* perform quantification as an integrated part of the inference process. Instead of producing an intermediate coordinate file (such as a BAM file), these tools directly estimate the number of fragments originating from each transcript. While the computational logic differs, the final output is functionally equivalent to traditional pipelines.

Regardless of the chosen method, the final output of NGS processing is the **raw count table**: a matrix where each row represents a biological feature (a gene or transcript) and each column represents an individual sample. The cell values are the discrete number of reads or fragments assigned to each feature. This raw count table is the fundamental output of NGS data processing and the starting point for all downstream statistical analyses

## Raw Counts Are Not Gene Expression Levels

Raw counts are *related to* gene expression levels, but they are not directly equivalent. Three technical factors influence raw counts independently of true expression:

- **Sequencing depth:** A sample sequenced to 30 million reads will show roughly twice the counts of the same sample sequenced to 15 million reads, for every gene. This does not reflect a biological difference.

- **Gene length:** Longer genes generate more fragments during library preparation and therefore accumulate more reads, even if they are expressed at the same level as shorter genes.

- **GC content:** Genes with extreme GC content may be amplified less efficiently during PCR, leading to systematic underrepresentation in the count table.

For comparing the same gene across samples (differential expression analysis), sequencing depth is the main concern — handled by normalization methods like *DESeq2*'s size factors or

*edgeR*'s TMM normalization. For comparing different genes within a sample, gene length must also be accounted for, which is why metrics like RPKM, FPKM, and TPM were developed. However, for the statistical analyses covered in this book, we work with raw counts and let the statistical models handle normalization internally.

> **[BEST PRACTICE: Start with Raw Counts]** Always begin your analysis with the raw count table, not with pre-normalized values like RPKM or TPM. Modern statistical frameworks (*DESeq2, edgeR, limma-voom*) are designed to work directly with raw counts and perform their own internal normalization.

## Key Practical Considerations

**Paired-end vs. single-end reads.** Paired-end sequencing reads each fragment from both ends, producing two reads per fragment. For gene expression quantification, single-end reads of 75–100 bases are typically sufficient — the k-mer analysis shows that 25–35 bases suffice for unique mapping in most genomes. Paired-end sequencing nearly doubles the cost with marginal benefit for standard expression analysis. However, paired-end reads are valuable for detecting structural variants, novel transcripts, and fusion genes.

**Sequencing depth.** For bulk RNA-seq in human samples, 10–20 million reads per sample captures the majority of the transcriptome, with diminishing returns beyond 25 million. Modern studies often sequence 30–40 million reads, which is more than adequate. The number of biological replicates matters more than sequencing depth: five or more replicates per group is recommended, as statistical power increases more rapidly with additional replicates than with additional sequencing depth.

**Multiplexing and demultiplexing.** Modern sequencing instruments are so powerful that a single run can generate billions of reads — far more than any single experiment requires. To make efficient use of this capacity, multiple samples are pooled (multiplexed) on a single run, each tagged with a unique **barcode** sequence. After sequencing, the reads are separated (demultiplexed) back to individual samples based on their barcodes. This is typically handled by the sequencing center, but it requires careful record-keeping of which barcode corresponds to which sample. One rule to follow: only multiplex samples of the same species together. Mixing human and bacterial samples on the same run creates amplification biases.

The following flowchart summarizes the complete NGS pipeline, from raw FASTQ files through quality control, alignment (or assembly), and feature counting, to the raw count table that serves as input for all downstream analyses.

*Figure 2.6 — NGS processing pipeline overview. A flowchart showing: Raw FASTQ files → Quality control (fastp) → Clean reads → Read mapping (HISAT2/STAR) or De novo assembly (SPAdes/Seq2Fun) → Read counting (featureCounts/HTSeq) → Raw count table (Genes × Samples × Counts).*

## 2.4 Mass Spectrometry Raw Data Processing

Mass spectrometry-based omics — including metabolomics, proteomics, and exposomics — measures the physical state of a system by analyzing ionized molecules. Unlike the discrete digital counts of NGS, MS data is analog and continuous, consisting of spectral peaks that must be distilled into biological features. Understanding this workflow is essential for interpreting the feature tables that emerge from this measurement engine.

### The LC-MS/MS Workflow: Separation, Ionization, and Identification

A liquid chromatography–tandem mass spectrometry (LC-MS/MS) experiment typically follows these three fundamental stages to transform a complex biological mixture into a multidimensional data array:

1. **Chromatographic separation (LC):** The sample is injected into a column (often a C18 column for proteomics, or diverse chemistries for metabolomics) to separate molecules in the time dimension. This prevents co-elution, where too many molecules hit the detector at once, and gives each analyte a characteristic **retention time (RT)**.

2. **Ionization:** Eluting molecules are converted into gas-phase ions, usually via **electrospray ionization (ESI)**. This assigns the molecules a mass-to-charge ratio (m/z). In proteomics, peptides often carry multiple charges (+2, +3), whereas metabolites typically carry a single charge (+1, −1).

3. **Mass detection (MS1 and MS2):** The instrument performs rapid, alternating cycles of two distinct scan types to capture both the quantity and the identity of the analytes

- **MS1 (precursor scan):** The instrument records the intensity and m/z of all intact precursor ions. In both fields, these data are the primary source for **quantification**. By integrating the signal intensity of a precursor ion across its elution time (the area under the peak), we determine its relative abundance in the sample

- **MS2 (fragmentation / product scan):** Specific precursor ions are isolated and energetically fragmented. The resulting fragment spectrum acts as a molecular "fingerprint" for **identification**. In proteomics, these fragments reveal the amino acid sequence of a peptide; in metabolomics, they represent the structural breakdown of a metabolite, which is matched against spectral libraries to confirm its chemical identity.

The three-stage LC-MS workflow — chromatographic separation, electrospray ionization, and mass detection — transforms a complex biological mixture into a structured three-dimensional data landscape of retention time, m/z, and intensity.

## Spectra Collection Strategy

The strategy for acquiring MS1 and MS2 data is dictated by the complexity of the molecular layer and the required depth of coverage. While both fields aim for comprehensive identification, they differ in how they prioritize instrument time across a study.

- **Shotgun Proteomics (Per-Sample MS2)**: In contrast, proteomics generally requires MS2 data for every sample to ensure confident protein identification. Traditionally, this was done via Data-Dependent Acquisition (DDA), where the instrument picks the top N most abundant peaks for fragmentation. However, modern proteomics has shifted toward Data-Independent Acquisition (DIA). Using high-speed instruments and deep-learning software, researchers systematically fragment all ions within wide mass windows across every sample. This approach significantly reduces the "missing value" problem and provides highly consistent, reproducible quantification across large cohorts

- **Metabolomics & Exposomics (Pooled MS2)**: A common strategy in large-scale studies is to perform MS1 (Full Scan) on every individual sample for quantification, while reserving MS2 (fragmentation) primarily for pooled quality control (QC) samples. These pools represent a representative mixture of the entire study. Within these pooled runs, researchers use either DDA or DIA to build a comprehensive spectral library. This library is then used to map identifications back to the features quantified in the individual MS1 runs. This strategy maximizes throughput and stabilizes the identification process across thousands of samples.

## Mass Accuracy and Resolution

In the MS-omics paradigm, the quality of biological discovery is governed by two fundamental properties of the mass analyzer: **mass accuracy** and **mass resolution**. These metrics dictate the reliability of the "distillation" process — transforming raw spectral peaks into confident biological identities.

### *Mass Accuracy: The Precision of Identity*

Mass accuracy defines how closely the measured m/z aligns with the theoretical value derived from a molecular formula (metabolomics) or a peptide sequence (proteomics).

- **In metabolomics:** Precision below 5 ppm is essential for determining the correct molecular formula of a small molecule from a database.

- **In proteomics:** High mass accuracy allows for the confident matching of fragmented peptides to a theoretical protein database. This precision is vital for distinguishing between peptides with nearly identical masses or identifying subtle post-translational modifications (PTMs).

> **[DEFINITION: Parts Per Million (ppm)]** The unit used to express mass accuracy, representing the relative deviation between the measured m/z and the theoretical m/z.
>
> $$\text{ppm} = \frac{(m/z_\text{measured} - m/z_\text{theoretical})}{m/z_\text{theoretical}} \times 10^6$$
>
> A lower ppm value (e.g., <1 ppm) allows researchers to distinguish between molecules with nearly identical masses, such as different peptide modifications or distinct **isobaric metabolites** (molecules with different chemical formulas but nearly identical masses). Modern high-resolution instruments, such as the **Orbitrap**, frequently achieve mass accuracy below 1 ppm. This mass accuracy provides the precise coordinate needed to anchor a signal to a specific chemical formula.

*Mass Resolution: Distinguishing Overlapping Signals*

Mass resolution describes the instrument's ability to distinguish two distinct peaks with very similar m/z values, typically measured as the full width at half maximum (FWHM) of a peak. Higher resolution results in narrower, sharper peaks. This allows the detector to resolve closely spaced ions that would otherwise merge into a single, uninterpretable signal — effectively "cleaning" the data at the point of acquisition.

- **In metabolomics:** High resolution is essential for distinguishing between isobaric compounds. It is also required to resolve complex adducts and fine isotopic patterns (e.g., distinguishing a $^{15}N$ isotope from a $^{13}C$ isotope), which is critical for formula prediction.
- **In proteomics:** High resolution is vital for resolving the isotopic envelope of a peptide. Because peptides are often multi-charged, the distance between their isotopes ($\Delta$ m/z) becomes very small. High resolution is necessary to clearly separate these isotopes to accurately assign the peptide's charge state and mass.

*Analyzer Comparison: TOF vs. Orbitrap*

- **Time-of-Flight (TOF):** Offers extremely fast acquisition speeds with moderate-to-high resolution. It is often used in high-throughput screening or when rapid duty cycles are required.
- **Orbitrap:** Provides the highest available resolution and mass accuracy at somewhat slower acquisition speeds. For untargeted omics — both metabolomics and proteomics — where comprehensive coverage and deep feature identification are the goals, Orbitrap analyzers are the preferred instrument.

The visual difference between low, moderate, and high mass resolution is striking: two closely spaced peaks that appear as a single broad hump at low resolution become fully separated and individually quantifiable at high resolution.

*Figure 2.8 — Mass resolution comparison. Three panels showing the same two closely spaced peaks at low resolution (merged into one broad peak), moderate resolution (partially resolved), and high resolution (fully resolved into two sharp, narrow peaks). Labels indicate FWHM and the importance of resolution for metabolomics.*

## 2.5 MS Data Processing and Feature Identification

Raw LC-MS data — a continuous three-dimensional landscape of **RT, m/z, and intensity** — is uninterpretable in its "analog" state. To perform statistical analysis, this data must be "distilled" into a discrete feature table. This process involves several critical computational steps.

1. **Peak Picking (Detection)** The first step is peak picking, where algorithms scan the RT–m/z–intensity space to identify local maxima and distinguish biological signals from background chemical noise. In metabolomics, each detected peak typically represents a specific ion or adduct of a small molecule. In proteomics, each peak represents an individual isotopic component of a peptide. The resulting output is a list of detected peaks, each characterized by its centroid m/z, elution time, and peak area (the integrated signal intensity), which serves as a proxy for abundance.

2. **Peak Alignment and Grouping** Because retention times naturally drift between runs due to factors like column aging or temperature fluctuations, peaks must be aligned across the entire sample set. Alignment algorithms correct for these shifts so that the same molecule across dozens of samples is correctly assigned to the same data row. Following alignment, the process of grouping (including de-isotoping and de-charging) occurs. In both fields, but particularly in proteomics, multiple peaks (isotopes) belong to the same parent molecule. Algorithms group these related signals into a single "feature" to prevent redundant counting and ensure that the final quantification reflects the total abundance of the molecule rather than just one of its isotopes.

3. **Feature Table Generation.** The final product is a matrix where each row is a unique molecular feature and each column is a sample. The final product of this pipeline is a matrix where each row represents a unique molecular feature and each column represents a sample. Each row is typically assigned a **Feature ID** based on its m/z and RT (e.g., 401.05_2.3 for m/z = 401.05 at 2.3 min). This ID is the MS-omics equivalent of a gene ID in an NGS count table, providing a unique identifier for statistical testing even if the specific chemical or protein identity remains unconfirmed at this stage.

The peak processing pipeline transforms a continuous three-dimensional spectral landscape into a discrete, analysis-ready feature table through peak picking, retention-time alignment, and grouping.

*Figure 2.9 — Peak processing pipeline. Panels showing: (1) Raw 3D data (RT × m/z × intensity landscape), (2) Peak picking identifying individual peaks, (3) Peak alignment correcting RT shifts across samples, (4) Final feature table with rows labeled as m/z_RT and columns as samples.*

## Feature Identification: Levels of Confidence

A peak in the feature table is not synonymous with an identified biological entity. In the MS-omics paradigm, moving from a numerical coordinate to a biological identity is a process of establishing confidence levels based on the strength of available evidence. While the general principles of spectral matching are shared across the "omes", the specific criteria for

metabolomics and proteomics are distinct and represent different tiers of chemical or biological certainty.

- **Level 1 — Confirmed Identification** represents the highest tier of certainty, where the identity is confirmed by matching experimental data against an authentic reference standard analyzed on the same instrument under identical conditions. In metabolomics, this requires a three-way match across accurate mass, MS/MS fragmentation, and retention time (RT), serving as the "gold standard" for chemical certainty. In proteomics, this level is typically reserved for targeted workflows, such as Parallel Reaction Monitoring (PRM) or Selected Reaction Monitoring (SRM), where a match to a synthesized peptide sequence is used to quantify specific biomarkers with absolute certainty.

- **Level 2 — Putative Identification** is the standard for discovery-based research, where identities are tentatively assigned by matching experimental data against large-scale digital databases rather than physical standards. Metabolomics researchers match accurate mass and MS/MS spectra against libraries like METLIN or HMDB. In shotgun proteomics, this involves matching fragment patterns against a theoretical database derived from the genome using search engines like *MaxQuant* or *DIA-NN*.

- **Level 3 — Putative Class or Protein Group** is assigned when only the general category of a molecule is clear. In metabolomics, a molecule might be identified as a "phospholipid" or "flavonoid" based on a characteristic chemical fragment, even if the exact structure remains unknown. In proteomics, this level often addresses the protein inference problem: a peptide is identified with high confidence, but because that sequence exists in multiple protein isoforms, it cannot be uniquely mapped to a single source, resulting in a "protein group" designation.

## The "Dark Matter" of Untargeted Omics

In untargeted studies, a significant portion of the data remains unannotated, representing the "dark matter" of the molecular landscape. In metabolomics and exposomics, it is common for only 10–30% of peaks to receive a putative identification; the remainder consists of real biological signals from molecules not yet characterized in our libraries. Identification rates in proteomics are generally higher (50–80%) due to the predictable nature of the genetic code, but many spectra remain unassigned due to unexpected post-translational modifications or non-canonical protein translations.

As explored in later chapters, a lack of 100% identification does not prevent meaningful biological discovery. Modern statistical enrichment methods are designed to be robust, looking for coordinated patterns across entire pathways or protein networks. This allows the collective

signal of identified features to overcome the noise of the unknowns, enabling researchers to distill biological insight even from incomplete data.

## Relative vs. Absolute Quantification

In the MS-omics paradigm, the intensities recorded in the feature table—whether measured as peak areas or peak heights—reflect relative abundance rather than absolute concentration. While a higher peak generally indicates a higher quantity of a molecule, two primary chemical phenomena prevent a direct, linear conversion to concentration across different features.

- **Differential Ionization Efficiency.** Molecules do not ionize equally. The chemical structure of a metabolite or the specific amino acid sequence of a peptide dictates how "easy" it is for that molecule to carry a charge. In metabolomics, a metabolite present at 1 μM may produce a significantly larger peak than another at 10 μM simply because it is more amenable to the electrospray process. Similarly, in proteomics, different peptides originating from the exact same protein can have vastly different intensities due to variations in their length, hydrophobicity, and charge-carrying capacity.

- **Matrix Effects and Ion Suppression.** The biological matrix, such as plasma, fecal water, or tissue lysate, acts as a complex chemical background that influences every measurement. A common manifestation of this is ion suppression: if multiple compounds elute from the chromatography column at the same time, they compete for the limited charge available during ionization. High-abundance molecules can effectively "suppress" the signal of their lower-abundance neighbors. Consequently, the same molecule at the same concentration can produce wildly different peak areas when measured in plasma versus urine because of the differing background chemical "noise".

## From Discovery to Validation

The way these challenges are managed depends on the goal of the study:

1. **The Discovery Paradigm (Untargeted)**: Much like NGS counts, untargeted MS-omics focuses on changes between groups. During statistical analysis, data are transformed and normalized to account for technical variation, making relative comparisons mathematically valid even without knowing absolute concentrations.

2. **The Validation Paradigm (Targeted)**: When absolute concentration is required (e.g., for clinical diagnostics), researchers use targeted methods such as SRM or PRM. These workflows utilize isotope-labeled internal standards—the "heavy" version of a molecule—which are spiked into the sample to provide an absolute, known reference point for quantification.

> **[BEST PRACTICE: Relative Abundance Is Sufficient for Discovery]** For the broad discovery phase of omics — where the goal is to compare relative differences between experimental groups (e.g., diseased vs. healthy) — absolute quantification is rarely necessary.

Despite these fundamentally different measurement technologies, NGS count tables and MS intensity tables share the same features-by-samples architecture. This structural similarity allows them to feed into a unified downstream pipeline for statistical analysis and functional interpretation, which we will explore in the next section.

*Figure 2.10 — Comparison of NGS and MS feature tables. Left: NGS count table (Gene ID × Samples × Raw Counts). Right: MS feature table (m/z_RT × Samples × Peak Intensity). Both are matrices of features by samples by abundance values, despite arising from fundamentally different measurement technologies. Arrows show that both feed into the same downstream statistical analysis pipeline.*

## 2.6 From Features to Functions: Enrichment Analysis

At this point in the data distillation process, raw data has been processed into a feature table, and statistical analysis has identified a set of significant features — genes whose expression differs between conditions, or metabolite that distinguish disease from control. But a list of 500 significantly changed genes is not, by itself, a biological insight. It remains at the feature level.

When **enrichment analysis** reveals that these genes converge on immune response, cell cycle regulation, and metabolic reprogramming, isolated data points coalesce into the cross-layer principles this book is built around. Functional interpretation is where the distillation process reaches its destination — transforming numbers into rich biological context.

### How Functional Analysis Works: The Power of Group Behavior

Functional analysis is fundamentally about **collective behavior**. Rather than interpreting individual genes or metabolites in isolation, we ask whether groups of functionally related molecules show coordinated changes. This shift from individual features to functional groups is another form of dimension reduction — from thousands of individual molecules to dozens of pathways or functional categories.

The strength of the omics paradigm lies in collective signal over individual precision. This group-level analysis has a remarkable property: it is robust to individual errors. Even if individual feature annotations contain mistakes, the collective signal of a pathway remains detectable as long as the errors are random rather than systematic. Simulation studies have

quantified this tolerance: with as little as 30% correct annotation, enrichment analysis can still detect 100% of the truly affected pathways, provided the errors are randomly distributed. Three reasons explain this robustness:

- **Tolerance for uncertainty:** Functional and pathway enrichment analyses are designed to identify coordinated shifts across entire biological modules (e.g., the TCA cycle or a protein signaling pathway).

- **The random vs. systematic rule:** As long as annotation errors are random (distributed across the data), they act as "background noise". The true biological "signal" — driven by the features that are correctly identified — remains statistically visible.

- **Patterns over particulars:** In the data distillation framework, our goal is to identify robust biological functions rather than perfecting the identity of every single coordinate in the feature table.

> **[TIP: Approximately Correct Is Often Sufficient]** The principle of functional robustness is a cornerstone of omics data science. While individual feature annotations can often be incomplete or inaccurate—such as in untargeted metabolomics, where confident identification rates may be as low as 10–30%—this does not prevent the discovery of reliable biological insights. This is because biological systems operate through highly interconnected networks and pathways, system-wide functional patterns remain robust and interpretable even when a fraction of the individual components are missing or misidentified. The aggregate signal of the pathway overrides the noise of individual annotation errors.

## Three Components of Enrichment Analysis

Every enrichment analysis requires three ingredients:

1. **Signal definition (input):** What features do you want to test? This can be a list of significant features (e.g., genes with $p < 0.05$), a ranked list of all features (e.g., genes ranked by fold-change), or features from a specific cluster or pattern.

2. **Functional library:** What biological categories do you want to test against? Common libraries include Gene Ontology (GO), KEGG pathways, Reactome, and metabolite set libraries. Each category contains a predefined list of genes or metabolites known to participate in a particular biological function.

3. **Algorithm:** How do you determine whether a category is enriched? The algorithm compares the overlap between your signal features and each functional category to what would be expected by chance.

## Functional Libraries

The choice of functional library determines what biological questions you can ask. The most commonly used libraries include:

- **Gene Ontology (GO):** A hierarchical classification of gene functions organized into three categories: Biological Process, Molecular Function, and Cellular Component. GO is the most comprehensive functional annotation for genes and proteins.

- **KEGG (Kyoto Encyclopedia of Genes and Genomes):** Curated maps of metabolic and signaling pathways, widely used for both genomics and metabolomics enrichment analysis. KEGG pathways are among the most frequently used resources in omics studies.

- **Reactome:** A curated database of human biological pathways and reactions, with extensive annotations for signaling, metabolism, and disease pathways. Reactome provides detailed mechanistic information not always captured in KEGG.

- **MSigDB (Molecular Signatures Database):** A collection of thousands of curated gene sets organized into categories including hallmark gene sets, positional gene sets, regulatory targets, and computational gene sets. Widely used with *GSEA*.

- **Metabolite sets:** Collections of metabolites grouped by biological function, disease association, chemical class, or localization. Available through tools like *MetaboAnalyst*.

- **Taxon sets:** Curated collections of microbial taxa grouped by their shared ecological roles, clinical relevance, or biological characteristics. Available through tools like *MicrobiomeAnalyst*.

The size of functional categories matters for statistical power. Categories that are too small (fewer than 5 members) lack power — even if all members are significant, the result may not be statistically convincing. Categories that are too large (more than 200 members) are too general to provide specific biological insight. The optimal range is approximately **15–80 members**, which provides both statistical power and biological specificity.

> **[BEST PRACTICE: Filter by Set Size]** When running enrichment analysis, filter out functional categories with fewer than 5 members (insufficient

power) or more than 200 members (too general). Most enrichment tools allow you to set minimum and maximum set size thresholds.

## 2.7 Enrichment Algorithms

Enrichment analysis is the computational bridge that transforms a feature table of thousands of measurements into a handful of biological stories. Before selecting an algorithm, one must choose a statistical "frame of reference":

- **Competitive Methods (Discovery Standard):** These algorithms ask, "Is this pathway more affected than the rest of the features in my experiment?" They use the global background of all measured molecules as a baseline. This is the gold standard for discovery science because it filters out global biological "noise" and prioritizes uniquely important themes.

- **Self-Contained Methods (Hypothesis Testing):** These algorithms ask, "Is there any significant change happening within this specific pathway?" They ignore the rest of the data. While mathematically more powerful, they can be overly sensitive in noisy datasets, potentially flagging nearly every pathway as significant.

### The "Urn" Model: Over-Representation Analysis (ORA)

ORA is the most intuitive competitive method. It treats your features like colored balls in an urn. If you draw a subset of "significant" features (the balls), are there more members of "Pathway A" in that draw than you would expect by random chance?

- **Logic:** Uses the hypergeometric test (or Fisher's exact test) to evaluate a 2 × 2 contingency table.
- **Input:** Requires a "hard" cutoff (e.g., $p < 0.05$ and |fold-change| > 2).
- **Best for:** A fast "sanity check" when you have a clear, high-confidence list of differentially expressed features.

### The "Trend" Model: GSEA

A major limitation of ORA is its dependence on arbitrary cutoffs, which can exclude genes that are biologically important but fall just below a $p$-value threshold. Gene Set Enrichment Analysis (GSEA) is a competitive method that avoids cutoffs by analyzing the entire ranked dataset. It asks: *"Do the members of this pathway generally trend toward the top or bottom of my ranked list?"*

- **Logic:** Calculates a running enrichment score (ES). If pathway members cluster at the extremes of your ranked list, the score peaks.

- **Input:** A complete ranked list of all measured features (e.g., ranked by t-statistic or fold-change).

- **Best for:** When biological changes are subtle or distributed across many features, ensuring you do not "lose" signal due to a strict *p*-value threshold.

### The "Pattern" Model: GlobalTest

*GlobalTest* represents a **self-contained** mathematical approach. Instead of looking at individual feature *p*-values or ranks, it examines the global expression pattern of a group of features across your samples to see if it correlates with your experimental groups.

- **Logic:** A regression-based approach that tests whether the collective expression pattern of a feature set is significantly associated with the clinical outcome or experimental group.

- **Input:** The normalized data matrix for all features within a set.

- **Best for:** When you have a specific hypothesis about a pathway and want maximum sensitivity to detect if *anything* in that pathway is changing in response to your condition.

> **[BEST PRACTICE: Choosing Your Enrichment Strategy]** The choice of an enrichment strategy depends on whether you are prioritizing **discovery** or **biological sensitivity**. Use **competitive** methods for initial, unbiased exploration where you need to filter thousands of features down to the most uniquely active themes. Shift to **self-contained** methods when you have a pre-defined biological hypothesis or when working with sparse datasets where measurement bias might lead to high false-negative rates.

In a *GSEA* analysis, the enrichment score curve rises steeply when pathway members cluster among the most upregulated genes and falls when they are absent, with the maximum deviation from zero quantifying the strength of coordinated pathway regulation.

*Figure 2.11 — GSEA illustration. A ranked list of all genes from most upregulated (left) to most downregulated (right). A specific pathway's members are marked with vertical lines. The running enrichment score (a curve) increases when a pathway member is encountered and decreases otherwise. The maximum deviation from zero is the enrichment score (ES). The pathway shown has members concentrated among the upregulated genes, producing a positive ES.*

## Potential Biases in Enrichment Analysis

Enrichment analysis is powerful but susceptible to several systemic biases that can lead to spurious biological conclusions.

- **Signal Selection Bias:** ORA results are highly sensitive to the significance threshold (e.g., $p < 0.05$ vs. $p < 0.01$) used to define the signal list. When a list is small, a single feature's presence or absence can cause entire pathways to appear or disappear.

- **Background Definition:** Enrichment is calculated relative to a **background**—the set of all features that *could* have been detected. Using an inappropriate background (e.g., the entire genome instead of only tissue-specific expressed genes) artificially inflates significance.

- **Systematic and Biological Bias:** The hypergeometric test assumes every feature has an equal probability of being detected or targeted. This is often violated:

    - **In MS-Omics:** Certain chemical classes are preferentially detected due to column chemistry or ionization efficiency.

    - **In miRNA Research:** The "standard" hypergeometric method is often inappropriate for miRNA target genes because certain functional categories are preferentially targeted by miRNAs in general, regardless of the specific biological condition. For example, genes with longer 3' UTRs are more likely to be predicted as miRNA targets, creating an underlying bias in the distribution.

- **Database Coverage Bias:** Algorithms can only "see" what is annotated. Technologies like proteomics typically detect only 25-30% of the known proteome, biasing results toward well-studied, high-abundance pathways.

> **[CAUTION: Enrichment Is Not Importance]** Enrichment $p$-values reflect statistical significance, not biological importance — the same principle behind the significance vs. performance distinction introduced in Chapter 1. A highly significant pathway may be an artifact of database structure or detection bias. Expert biological judgment and literature validation remain essential for interpreting enrichment results.

## Addressing Bias: Permutation and Empirical Approaches

To overcome these biases in enrichment tests, researchers increasingly use permutation and empirical sampling that build a null distribution directly from the experimental data.

- **Mummichog (Metabolomics):** In untargeted metabolomics, the mummichog algorithm bypasses the need for accurate compound identification. It maps significant m/z features directly to metabolic networks and uses permutation to assess if the observed local "activity" in a pathway is higher than what would happen by chance, effectively handling ambiguous identities and the high rate of "dark matter".

- **The Empirical Sampling Approach (miRNA):** To correct for the inherent targeting bias in miRNAs, researchers use an algorithm that randomly samples sets of miRNAs (rather than genes) of the same size as the experimental list. By comparing the observed enrichment to this empirical distribution, researchers find that the vast majority of "significant" categories reported in literature disappear after bias correction.

*The Permutation Procedure:*

1. Compute enrichment based on the original signal list.
2. Randomly draw a new list of the same size from the **measured features** (the actual background).
3. Compute enrichment for this random list.
4. Repeat this 1,000+ times to build an **empirical null distribution**.
5. The p-value is the proportion of times random draws outperformed your real data. If fewer than 5% of random permutations produce enrichment as strong as the original, the pathway is considered significantly enriched ($p < 0.05$).

Permutation approaches are particularly valuable when you want to test enrichment for any patterns — not just statistically significant features, but also features from a specific cluster, features with a particular expression trend, or features selected by visual inspection. Because the null distribution is derived empirically from your own data, it naturally accounts for the biases present in your specific experiment.

## Summary

This chapter has traced the full arc of omics data analysis, from the raw instrument output at the top of the distillation funnel to the biological functions at the bottom.

- Omics technologies fall into distinct data paradigms — NGS-omics (digital counts) and MS-omics (continuous intensities), with their recent extensions to spatial omics (coordinate data), and multi-modal omics (integrated paradigms) — each requiring specialized processing pipelines. Despite these different starting points, they all converge on a common feature table.

- **NGS raw data processing** transforms FASTQ files into gene count tables through quality control, read mapping (or *de novo* assembly), and feature quantification. While traditional workflows align reads to specific genomic coordinates, alignment-free methods use k-mer compatibility to skip the coordinate-mapping step entirely. Raw counts are influenced by sequencing depth, gene length, and GC content; they are digital proxies for abundance, but are not yet equivalent to absolute expression levels.

- **MS raw data processing** transforms continuous spectra into discrete feature tables through chromatographic separation, ionization, mass detection, peak picking, and alignment. The relationship between detected peaks and biological entities is many-to-many: one metabolite often produces multiple ion adducts, and one protein is represented by multiple unique peptides.

- Both NGS counts and MS intensities reflect **relative abundance** rather than absolute concentration. While these values are ideal for comparing relative differences between groups, they do not directly represent molarity or exact cellular concentrations without the use of internal standards or specific calibration methods.

- **Enrichment analysis** connects significant features to biological functions through strategies such as *Over-Representation Analysis*, *Gene Set Enrichment Analysis*, or *GlobalTest*. These methods allow researchers to transition from a list of individual features to a cohesive biological narrative, though they require careful attention to signal definition, background choice, and potential technical biases.

With the raw data processing pipeline (this chapter) and the normalization and statistical analysis framework (Chapter 1) in hand, you now have a complete conceptual understanding of the four-step omics workflow.

• • • • •

**Table 2.3 — Key Terms**

| Term | Definition |
|---|---|
| **FASTQ** | File format storing both nucleotide sequences and per-base quality scores from NGS instruments. |
| **Phred Score (Q score)** | Quality metric expressing the probability of a base call error on a logarithmic scale (Q30 = 99.9% accuracy). |
| **K-mer** | A contiguous subsequence of length *k* extracted from a read; used for indexing, mapping, and assembly. |
| **Read Mapping** | Aligning sequencing reads to a reference genome to determine their genomic origin and count gene expression. |
| ***De Novo* Assembly** | Reconstructing a genome or transcriptome from sequencing reads without a reference, typically using k-mer-based algorithms. |
| **Retention Time (RT)** | The time a compound takes to pass through a chromatographic column; characteristic of compound–column interaction. |
| **Mass-to-Charge Ratio (m/z)** | The ratio of an ion's mass to its charge; the fundamental measurement in mass spectrometry. |
| **Adduct** | An ion formed when a molecule associates with another ion ($Na^+$, $K^+$, $NH_4^+$) during electrospray ionization. |
| **Peak Picking** | The process of identifying discrete peaks in continuous LC-MS data, each representing a detected ion. |
| **Over-Representation Analysis (ORA)** | Enrichment method testing whether significant features are disproportionately represented in functional categories using the hypergeometric test. |
| **Gene Set Enrichment Analysis (GSEA)** | Enrichment method using a complete ranked list of features and a running sum statistic to detect coordinated changes in gene sets. |
| **Self-Contained Methods** | Enrichment methods that test whether features within a specific gene set show significant effects, without comparing to a background (e.g., GlobalTest). |
| **Competitive Methods** | Enrichment methods (ORA, GSEA) that compare a gene set against the background of all measured features to identify uniquely affected pathways. |

| Term | Definition |
|---|---|
| **Enrichment Score** | The maximum deviation of the GSEA running sum from zero; measures how strongly a gene set is concentrated at the extremes of a ranked list. |
| **Electrospray Ionization (ESI)** | A soft ionization technique that adds or removes protons to create charged ions from molecules eluting from a chromatographic column. |

• • • • •

## Review Questions

**1.** A colleague argues that single-end sequencing at 75 bp is wasteful because the reads are "too short" for accurate gene expression quantification. Using the k-mer concept, explain why 75 bp reads are more than sufficient for mapping to the human genome.

**2.** You receive a metabolomics feature table where features are labeled as m/z_RT values (e.g., "401.05_2.3"). Explain why a single compound might appear as multiple rows in this table, and why a single row might represent contributions from multiple compounds.

**3.** Your RNA-seq experiment produced 40 million paired-end reads per sample. A reviewer suggests that the raw count of Gene X (which is 3 kb long) cannot be directly compared to the raw count of Gene Y (which is 0.5 kb long) to determine which is more highly expressed. Explain why the reviewer is correct and what factors influence raw counts beyond true expression level.

**4.** You run ORA on a list of 50 significant metabolites and find that "amino acid metabolism" is the top enriched pathway. When you change your significance cutoff from $p < 0.05$ to $p < 0.01$, the signal list shrinks to 15 metabolites and "amino acid metabolism" is no longer significant. What does this instability suggest, and how might a permutation-based or *GSEA* approach help?

**5.** A mass spectrometry experiment using a C18 reverse-phase column consistently shows "lipid metabolism" as the top enriched pathway across unrelated experiments. Could this represent a systematic detection bias rather than a true biological finding? Explain.

**6.** Why does the *GSEA* approach not require setting a significance threshold for defining the input gene list, and in what situations might this be advantageous compared to ORA?

• • • • •

## Further Reading

- Conesa, A., *et al.* "A survey of best practices for RNA-seq data analysis." *Genome Biology* 17, 13 (2016).

- Pang, Z., *et al.* "*MetaboAnalyst* 6.0: narrowing the gap between raw spectra and functional insights." *Nucleic Acids Research* 52, W1 (2024).

- Subramanian, A., *et al.* "Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles." *Proceedings of the National Academy of Sciences* 102, 43 (2005): 15545–15550.

- Lu, Y., Pang, Z., and Xia, J. "Comprehensive investigation of pathway enrichment methods for functional interpretation of LC–MS global metabolomics data", *Briefings in Bioinformatics*, 24, 1 (2023): bbac553.

- Xia, J., and Wishart, D.S. "MSEA: a web-based tool to identify biologically meaningful patterns in quantitative metabolomic data." *Nucleic Acids Research* 38, suppl_2 (2010): W71–W77.

- Xia, J., and Wishart, D.S. "MetPA: a web-based metabolomics tool for pathway analysis and visualization." *Bioinformatics* 26, 18 (2010): 2342–2344.

- Li, S., *et al.* "Predicting network activity from high throughput metabolomics." *PLoS Computational Biology* 9, 7 (2013): e1003123.

- Bolger, A.M., Lohse, M., and Usadel, B. "*Trimmomatic*: a flexible trimmer for Illumina sequence data." *Bioinformatics* 30, 15 (2014): 2114–2120.

- Dobin, A., *et al.* "*STAR*: ultrafast universal RNA-seq aligner." *Bioinformatics* 29, 1 (2013): 15–21.

- Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. "Graph-based genome alignment and genotyping with *HISAT2* and HISAT-genotype." *Nature Biotechnology* 37, 8 (2019): 907–915.

- Prjibelski, A., *et al.* "Using *SPAdes de novo* assembler." *Current Protocols in Bioinformatics* 70, 1 (2020): e102.

## What's Next

We are now ready to dive into individual omics disciplines. Chapter 3 begins with transcriptomics — the most mature and widely practiced form of omics analysis — covering the evolution from microarrays to RNA-seq, differential expression workflows, and the practical tools for turning raw count tables into biological discoveries.