

# A Deep Dive into Single-Cell V(D)J Sequencing Technology: Background, Cutting-Edge Computational Methods, and Public Dataset Resources

## I. Introduction: The Necessity of High-Resolution Immune Repertoire Analysis

### A. The Specificity Engine of the Adaptive Immune System: V(D)J Recombination

The random recombination of V(D)J genes is the molecular basis for the adaptive immune system's generation of a vast receptor diversity, making it possible to recognize nearly all antigens. The adaptive immune system is the primary defense mechanism in vertebrates for combating pathogens and eliminating abnormal cells, such as tumors. Its core features, high specificity and memory, are mainly mediated by T-cells and B-cells. These cells recognize antigens through their surface receptors: the T-cell receptor (TCR) and the B-cell receptor (BCR), which is the membrane-bound form of an antibody. V(D)J gene recombination is the molecular basis for the astonishing diversity of these receptors. During T-cell and B-cell development, the V (Variable), D (Diversity, only in TCR $\beta$  and BCR heavy chains), and J (Joining) gene segments that encode the variable regions of TCRs and BCRs are randomly combined. Additionally, non-templated nucleotides (N/P nucleotides) are inserted at the junctions, generating a theoretically astronomical number of unique receptor sequences. It is this immense diversity that endows the immune system with the ability to recognize a virtually infinite array of antigens, allowing it to respond to ever-emerging pathogens and malignancies. Each T-cell or B-cell typically expresses a single, unique TCR or BCR, and this unique receptor sequence defines the cell's clonotype.

### B. Limitations of Traditional Bulk Repertoire Sequencing

Traditional bulk sequencing methods cannot resolve the native pairing of receptor chains, nor can they link sequences to the function of individual cells, severely limiting our in-depth understanding of immune responses. Before the advent of single-cell technologies, the study of immune repertoires primarily relied on bulk sequencing methods. These methods, which involve amplifying and sequencing the TCR or BCR genes from a large population of T- or B-cells, can provide information about the overall composition of the repertoire, such as V(D)J gene usage frequencies, CDR3 (Complementarity-Determining Region 3, the key region for antigen recognition) length distribution, and clonotype abundance. However, bulk sequencing has inherent limitations. First, it cannot resolve the native pairing of TCR  $\alpha$  and  $\beta$  chains (or BCR heavy and light chains) within a single cell. Since the function of a TCR depends on the correct  $\alpha$ - $\beta$  combination and that of a BCR on the heavy-light chain pair, the loss of this pairing information severely hampers an accurate understanding of antigen recognition specificity. Second, bulk sequencing mixes receptor sequences from different cells, making it impossible to directly link a specific TCR or BCR sequence to the phenotype or functional state of its cell of origin (e.g., activation state, differentiation subset, cytokine expression profile). This loss of information makes it extremely difficult to understand the specific role of a particular clone in an immune response.

### C. The Rise and Significance of Single-Cell V(D)J Sequencing

Single-cell V(D)J sequencing revolutionizes these limitations by simultaneously capturing paired receptor sequences and transcriptome information at the single-cell level, enabling unprecedented high-resolution analysis of immune responses. The emergence of single-cell V(D)J sequencing technology represents a revolutionary breakthrough in overcoming the limitations of traditional bulk sequencing. This technology enables the simultaneous acquisition of full-length, paired TCR  $\alpha$  and  $\beta$  chain (or BCR heavy and light chain) coding sequences at the single-cell level. Furthermore, it can be combined with transcriptome sequencing (scRNA-seq) or cell surface protein sequencing (e.g., CITE-seq), thereby linking immune receptor sequence information with the gene expression profile, surface protein expression profile, and other molecular features of the same cell. This multi-modal, high-resolution analysis capability allows researchers to dissect the complexity of immune responses with unprecedented depth. For instance, it is possible to precisely identify T- or B-cell clones that undergo clonal expansion in specific disease states (like cancer or infection) and simultaneously understand their functional states (e.g., effector, memory, exhausted) and phenotypic characteristics. As some studies emphasize, obtaining "full-length, paired V(D)J sequences from B cells or T cells" and combining them with other data modalities is the key advantage of this technology.

The paradigm shift from population average to single-cell precision enables the study of rare clones and receptor function, with the integration of V(D)J sequence and transcriptome being its core value. This shift from population average to individual precision represents a fundamental paradigm shift in immune repertoire research. It is no longer just about making statistical inferences about the repertoire but about direct, high-resolution measurement of individual immune cells and their antigen receptors. This increase in granularity makes it possible to study rare clones, subtle phenotypic differences associated with specific receptors, and the true pairing of receptor chains, all of which were extremely challenging in the past. Therefore, single-cell V(D)J sequencing opens new avenues for a deeper understanding of the mechanisms of immune responses. It is crucial to emphasize that the ability to integrate V(D)J sequences with transcriptomic data at the single-cell level is not just an add-on feature but the core value driver of this technology. A TCR sequence by itself may reveal its potential antigen specificity, but by combining it with the T-cell's transcriptome, we can determine whether it is an effector, memory, exhausted, or regulatory T-cell. This combined information is vital for understanding key scientific questions, such as why certain T-cell clones expand during cancer immunotherapy while others do not, or how T-cell states evolve during an infection.

### D. Core Objectives of Single-Cell Immune Repertoire Research

Leveraging the power of single-cell V(D)J sequencing, research objectives typically revolve around five key areas: clonotype identification, clonality and diversity assessment, V(D)J gene usage analysis, and linking sequence to function. These objectives collectively form the foundational framework for exploring the complexity of the immune system using single-cell V(D)J sequencing.

- Clonotype Identification:** Precisely identifying the TCR or BCR sequence of each single cell to define its clonotype.
- Clonality Assessment:** Analyzing the relative abundance and distribution of different clonotypes in the immune repertoire to assess the degree of clonal expansion.
- Diversity Analysis:** Measuring the number and evenness of unique clonotypes in the immune repertoire to evaluate its breadth.
- V(D)J Gene Usage Patterns:** Studying the usage frequencies of different V, D, and J gene segments and their combinatorial preferences.
- Linking Sequence to Phenotype/Function:** Associating clonotype information with the cell's transcriptome, surface protein expression, or other functional properties to uncover the functional significance of specific clones.

## II. "Why We Study It": Rationale and Core Scientific Questions

The fundamental motivation for applying single-cell V(D)J sequencing is to decode the complex mechanisms of the adaptive immune response and to apply this understanding to various aspects of human health and disease.

### A. Foundational Goal: Understanding Antigen Specificity and Immune Memory

The core motivation for single-cell V(D)J sequencing is to decode the mechanisms of antigen-specific recognition and immune memory formation by tracking clonal dynamics to understand the fundamental operations of the immune system. The core of adaptive immunity lies in its specific recognition of antigens and the generation of lasting immunological memory. T-cells recognize antigenic peptides presented by the major histocompatibility complex (MHC) via their TCRs, while B-cells recognize native antigens directly through their BCRs. Single-cell V(D)J sequencing helps us to:

- Track Antigen-Specific Clones:** Following antigen exposure (e.g., infection or vaccination), identify and track the dynamic changes of T- and B-cell clones that specifically recognize the antigen, including their expansion, contraction, and differentiation into effector or memory cells.
- Investigate the Formation and Maintenance of Immune Memory:** By analyzing the clonotypic composition and phenotypic characteristics of memory T- and B-cells, we can uncover how immunological memory is established, maintained, and reactivated upon re-encountering the same antigen.
- Explore TCR/BCR Cross-Reactivity:** Some TCRs or BCRs can recognize multiple structurally similar but distinct antigens. This cross-reactivity plays an important role in both immune protection (e.g., against viral variants) and autoimmunity (e.g., misidentifying self-antigens). Single-cell analysis facilitates the study of this phenomenon at the clonal level.

Combining receptor sequences (the "vocabulary") with cellular function (the "speaker's characteristics") allows us to comprehensively interpret the precise role of specific immune clones in a biological context. At a deeper level, the language of adaptive immunity is composed of countless unique TCRs and BCRs (the "vocabulary") and the antigen-recognition specificities they represent (the "meaning"). Single-cell V(D)J sequencing, especially when combined with information on cell phenotype and function (the "speaker's" characteristics), allows us to more comprehensively interpret the precise meaning and role of this "vocabulary" in specific immune contexts. For example, identifying a tumor-reactive TCR sequence is the first step; understanding whether the T-cell expressing this TCR is a killer effector cell or a functionally exhausted cell, by integrating transcriptomic data, is the crucial second step to understanding its true role in anti-tumor immunity. This integrated analysis is essential for elucidating disease mechanisms or evaluating treatment efficacy.

### B. Applications in Disease Contexts

Single-cell V(D)J sequencing shows immense potential in the study of various diseases, including oncology, infectious diseases, and autoimmune disorders.

#### 1. Oncology

In oncology, this technology is used to identify and characterize key anti-tumor immune cells, track responses to immunotherapy, and provide candidate receptors for developing novel cell therapies.

- Identifying Tumor-Infiltrating Lymphocytes (TILs):** TILs are key effector cells in the anti-tumor immune response. Single-cell V(D)J sequencing can identify the TCR clonotypes in TILs, reveal the extent of T-cell clonal expansion in the tumor microenvironment, and analyze their functional states (e.g., effector, exhausted, memory) by integrating transcriptomic data.
- Tracking Response to Immunotherapy:** During immunotherapy, such as with immune checkpoint inhibitors, tracking the dynamic changes of T- and B-cell clones in patients can help to understand the mechanisms of response, distinguish responders from non-responders, and potentially discover new therapeutic targets or biomarkers.
- Developing TCR-Engineered T-cell (TCR-T) Therapies:** By identifying TCRs with high affinity and specificity for tumor antigens, these TCRs can be used to engineer T-cells from patients or healthy donors for adoptive cell therapy.

#### 2. Infectious Diseases

In infectious disease research, single-cell analysis helps us identify T-cells that can effectively clear pathogens and B-cells that produce neutralizing antibodies, as well as evaluate vaccine efficacy.

- Characterizing Pathogen-Specific Immune Repertoires:** During viral (e.g., HIV, influenza, SARS-CoV-2) or bacterial infections, analyzing the changes in TCR and BCR repertoires can identify clonotypes targeting specific pathogen antigens.
- Understanding Protective Immune Responses:** Studying the characteristics of neutralizing antibodies (BCR sequences) and effector T-cells (TCR sequences) that are associated with protective immunity in recovered patients or vaccinated individuals.
- Monitoring Vaccine Responses:** Assessing the breadth, magnitude, and durability of the immune response induced by a vaccine, and identifying key clonotypes associated with protective efficacy.

3. Autoimmune Diseases

For autoimmune diseases, this technology is a powerful tool for identifying and understanding the culprits—autoreactive T/B cells—that mistakenly attack self-tissues.

- **Identifying Autoreactive T- and B-cell Clones:** In autoimmune diseases such as rheumatoid arthritis, systemic lupus erythematosus, and type 1 diabetes, identifying the T- and B-cell clones that mistakenly attack self-tissues.
- **Understanding Mechanisms of Broken Immune Tolerance:** By analyzing the characteristics of autoreactive clones and their regulatory environment, we can investigate how immune tolerance is broken, leading to autoimmunity.
- **Tracking Therapeutic Efficacy:** Evaluating the effectiveness of treatments aimed at depleting or modulating autoreactive clones.

C. Vaccine Development and Efficacy Assessment

In the field of vaccine development, single-cell V(D)J sequencing can precisely assess the quality of vaccine-induced immune responses, thereby guiding and optimizing the design of next-generation vaccines.

- Assess the breadth and depth of the immune response by analyzing the diversity, clonal expansion, and V(D)J gene usage patterns of the TCR and BCR repertoires induced by a vaccine.
- Identify clonotypes associated with protective immunity by comparing the immune repertoires of protected versus unprotected individuals after vaccination.
- Guide next-generation vaccine design by leveraging an understanding of the features of an effective immune response to engineer novel vaccines that more efficiently induce these desired characteristics.

D. Examples of Core Research Questions

Using single-cell V(D)J sequencing, researchers aim to answer a series of key scientific questions, such as the characteristics of specific antigen receptors, the dynamic changes of the immune repertoire in disease, and the existence and significance of public clones.

- What are the sequence and structural features of TCRs/BCRs that recognize specific antigens (e.g., tumor neoantigens, viral peptides)?
- How do the diversity and clonality of the immune repertoire change during disease progression or in response to treatment?
- Are there "public" TCRs/BCRs (i.e., clonotypes shared among different individuals) associated with specific diseases or immune responses?

The discovery of such public clones is highly significant for developing universal diagnostic tools or "off-the-shelf" cell therapies. If multiple individuals utilize similar or identical TCRs/BCRs to combat the same pathogen or cancer, these receptors are likely targeting key, conserved epitopes. This makes them ideal candidates for developing TCR-T therapies, designing vaccines aimed at inducing these specific responses, or creating diagnostic tools to assess immune competence. Public databases play a crucial role in such studies.

- What are the differences in V(D)J gene usage patterns between healthy individuals and patients, or between different disease states?

The study of V(D)J gene usage, while seemingly descriptive, can reveal selection pressures or preferences imposed by genetic susceptibility, chronic antigen exposure, or the disease process itself, providing clues to the fundamental rules of repertoire formation and selection. Humans have a finite set of V, D, and J genes. Preferential usage or pairing under specific conditions (e.g., specific HLA types, chronic infections) indicates the presence of selective pressure or inherent biases in the recombination machinery or clonal selection. This helps us understand susceptibility to disease or response to vaccines.

- What are the phenotypic and functional characteristics of cells expressing a specific target TCR/BCR?

III. Cutting-Edge Computational Methodology for Single-Cell V(D)J Repertoire Analysis

Extracting biological insights from raw sequencing data relies on a series of complex bioinformatics analysis pipelines and advanced computational methods.

A. Overview of the Bioinformatics Analysis Pipeline

A typical single-cell V(D)J analysis workflow comprises six major steps: data processing, sequence alignment, clonotype definition, metric calculation, visualization, and multi-modal integration.

1. **Raw Data Processing:** Including demultiplexing and quality control.
2. **V(D)J Gene Segment Alignment and Sequence Assembly:** Aligning reads to reference V, D, and J gene segments and assembling full-length TCR/BCR variable region sequences.
3. **Clonotype Definition and Grouping:** Clustering cells with identical or similar receptors into clonotypes based on sequence features (primarily the CDR3).
4. **Immune Repertoire Metric Calculation:** Quantifying various features of the repertoire, such as clonality, diversity, and V(D)J gene usage.
5. **Visualization:** Displaying the structure and dynamics of the repertoire through various plots.
6. **Integration with Other Single-Cell Modalities:** Correlating V(D)J information with data from transcriptomics, surface proteomics, etc.

B. Raw Data Pre-processing and Quality Control (QC)

Data pre-processing and QC are the cornerstone of accurate downstream analysis, with the core tasks of correcting technical biases and accurately assigning sequencing reads to single cells. This step is crucial for ensuring the accuracy of subsequent analyses.

- **Cell Barcode Demultiplexing:** Assigning reads from a mixed sequencing pool to their single-cell of origin based on their cell barcodes.
- **Unique Molecular Identifier (UMI) Processing:** Using UMIs to correct for PCR amplification bias, allowing for more accurate quantification of the original number of each TCR/BCR transcript.
- **Filtering Low-Quality Reads and Cells:** Removing reads with low sequencing quality and "cells" with ambiguous barcodes or very low UMI counts to reduce noise.
- **Platform-Specific Initial Processing:** For example, the widely used 10x Genomics platform provides the Cell Ranger software suite, whose `cellranger vdj` command performs initial read processing, cell barcode assignment, UMI counting, and V(D)J sequence assembly and preliminary annotation.

C. V(D)J Gene Segment Annotation and Sequence Assembly

Accurate identification of V, D, and J gene segments and the hypervariable regions (especially CDR3) is a core step.

1. Alignment to Germline V(D)J Gene Segments:

Annotation of V(D)J gene segments primarily relies on aligning sequencing reads against a reference database of known germline genes. Classic tools like IMGT/HighV-QUEST, MiXCR, and IgBLAST are fundamental for identifying V, D, and J segments and extracting CDR3 sequences. They rely on alignment to databases of known germline genes. Some newer tools aim to improve alignment accuracy and speed, such as SONG (Systematic Optimization of Next-generation sequencing for Gapped-alignment), which is reported to have high accuracy and efficiency in V(D)J gene alignment.

2. De Novo Assembly of Full-Length Receptor Sequences:

For non-targeted scRNA-seq data, tools like TRUST4 can be used to assemble TCR/BCR sequences de novo, allowing for the acquisition of immune repertoire information even without V(D)J enrichment. For non-targeted single-cell RNA-seq (scRNA-seq) data, some tools can assemble TCR/BCR transcript sequences de novo from the total RNA sequencing reads. For example, TRUST4 can be used for de novo assembly from both bulk RNA-seq and scRNA-seq data. TraCeR and Platypus are specifically designed to reconstruct TCR sequences from scRNA-seq data, with Platypus paying special attention to cases of low TCR mRNA capture efficiency. 10x Genomics' Cell Ranger also performs sequence assembly.

D. Clonotype Definition and Grouping

A clonotype refers to a population of cells derived from a common lymphocyte progenitor and expressing the same TCR or BCR.

1. Defining a Clonotype:

Clonotype definition is typically based on shared CDR3 sequences, but the lack of a unified standard is a major challenge in the field, hindering direct comparison between studies. The most common definition is based on a shared CDR3 sequence. For paired TCR/BCR data, it is usually required that the CDR3 sequences of both the  $\alpha$  and  $\beta$  chains (or heavy and light chains) are identical. Sometimes, the requirement of identical V(J) gene usage is also added.

The AIRR (Adaptive Immune Receptor Repertoire) Community is actively promoting the development of data standards for immune repertoires, including the standardization of clonotype definition. This effort is crucial for facilitating data sharing and meta-analysis. Without a uniform standard, if Study A defines clonotypes based on CDR3 amino acid sequence and Study B defines them based on CDR3 nucleotide sequence and identical V gene usage, the reported number of clonotypes and level of clonal sharing could be vastly different, even when analyzing the same raw data. This inconsistency hinders efforts to build large-scale, integrable databases of clonotypes and their associated phenotypes/disease states. Standardization is key to ensuring reproducibility and fostering collaborative discovery.

2. Clonotype Identification Tools:

Various computational tools, such as scirpy and Dandelion, provide clonotype identification functions and can be integrated with mainstream single-cell analysis frameworks. Cell Ranger provides clonotype identification functionality within its analysis pipeline. scirpy offers flexible clonotype definition options (e.g., identical CDR3 nucleotide/amino acid sequence, or network clustering based on sequence similarity) and integrates tightly with scRNA-seq analysis workflows. Dandelion also performs clonotype identification and is designed to integrate with the popular scRNA-seq analysis package scanpy. VDJtools can process the output of various upstream alignment tools for clonotype identification and downstream statistical analysis. The Immcantation framework (including tools like Change-O, aLakazam) is highly powerful, especially in BCR analysis (e.g., somatic hypermutation analysis), and many of its concepts and methods can also be applied to TCR analysis.

E. Quantification of Immune Repertoire Metrics

To describe the immune repertoire at a macroscopic level, a series of quantitative metrics such as clonality, diversity, and V(D)J gene usage frequency need to be calculated. After identifying and grouping clonotypes, a series of metrics are computed to characterize the repertoire.

1. **Clonality:** Measures the skewness of the clonotype frequency distribution. High clonality usually indicates that the repertoire is dominated by a few highly expanded clones.
2. **Diversity:** Measures the richness (number) and evenness (frequency distribution) of unique clonotypes. Common indices include the Chao1 estimator, Shannon entropy, and Simpson's index.
3. **V(D)J Gene Usage:** Statistics on the frequency of different V, D, and J gene segments and their combinations.
4. **CDR3 Length Distribution:** Analysis of the length distribution of the hypervariable CDR3 region.
5. **Convergent Evolution:** Identifying identical or highly similar TCR/BCR sequences that arise independently, suggesting convergent selection for a common antigen.

Tools such as `scirpy`, `Dandelion`, `VDJtools`, and `alakazam` (part of the Immcantation suite) provide functions for calculating these metrics.

F. Immune Repertoire Visualization Techniques

Effective visualization techniques, such as clone frequency plots and network graphs, are key to intuitively understanding the structure and dynamics of complex immune repertoire data.

- **Clonotype Frequency Plots:** Bar plots or pie charts showing the most abundant clonotypes.
- **V(D)J Gene Usage Plots:** Heatmaps or chord diagrams displaying usage frequencies and combinations.
- **Diversity Rarefaction Curves:** Assessing the impact of sequencing depth on observed diversity.
- **Network Graphs:** Visualizing relationships between clonotypes based on sequence similarity (e.g., using `scirpy` and `Dandelion`).
- **Interactive Browsers:** E.g., the 10x Genomics Loupe V(D)J Browser for interactive exploration.

G. Integration with Other Single-Cell Modalities (e.g., scRNA-seq, CITE-seq)

The core advantage of single-cell V(D)J analysis lies in its ability to integrate receptor sequence information with other single-cell modalities like the transcriptome, thereby directly linking clonotype to the cell's biological state. This is the primary strength and a major direction of development in current single-cell immune repertoire analysis.

With the popularization of multi-omics technologies, analysis tools have evolved from standalone sequence analyzers to integrated solutions compatible with frameworks like Scanpy; this integration capability is fundamental for meaningful biological interpretation. The evolution of computational tools clearly reflects advancements in sequencing technology. As platforms like 10x Genomics made combined single-cell V(D)J and transcriptome sequencing more accessible, analysis tools evolved from standalone analyzers to integrated solutions (e.g., `scirpy`'s integration with `scanpy`). This integration is now standard for meaningful biological interpretation. The advent of single-cell multi-omics created an urgent need to link V(D)J information with other molecular data (e.g., gene expression), driving the development of new tools that can handle objects like `scanpy`'s `AnnData` and add V(D)J-derived metadata. Without this, researchers would be left with siloed datasets, losing the primary advantage of the technology.

1. Linking Clonotypes with Gene Expression:

Tools like `scirpy` and `Dandelion` are designed for tight integration with popular scRNA-seq analysis frameworks like `scanpy`. Frameworks like `Seurat` can also achieve this via custom scripts. This allows for overlaying clonotype information (e.g., clonal expansion) onto UMAP/t-SNE plots to visualize which cell subsets are enriched for specific clonotypes.

2. Cell Type Annotation:

Using gene expression data, tools like `CellTypist` or `SingleR` can annotate cell types (e.g., CD4+ naive T-cells). Repertoire characteristics can then be analyzed within specific subsets.

3. Integrating Cell Surface Protein Data (CITE-seq):

CITE-seq allows simultaneous measurement of RNA and surface protein expression. Combining V(D)J, transcriptome, and proteome data enables a more refined and accurate definition of cell states.

**Challenge:** Integrating multi-omics data from different sources requires careful handling of batch effects, effective normalization, and accurate alignment of information across modalities.

As datasets grow in size and complexity, the demand for scalable, robust, and user-friendly computational workflows becomes increasingly prominent. Well-documented, easy-to-install tools that can efficiently handle large-scale data will gain wider adoption. This also relates to the need for better visualization tools.

Simulation tools like `immuneSIM` are invaluable for benchmarking the accuracy and performance of new analysis methods, providing a controlled environment for methodological validation.

Given the inherent complexity and often unknown "ground truth" of real biological samples, simulated data provides a controlled environment for validation. Simulation tools allow us to generate repertoires with known properties (e.g., clone count, CDR3 sequences, error rates) and test how well different analysis tools can recover these known features. This is crucial for improving the reliability of computational methodologies.

Table 1: Comparative Overview of Single-Cell V(D)J Analysis Tools

Note: This table is an exemplary overview. Features may change with tool updates. Users should consult the latest documentation for detailed information.

Tool Name	Primary Function(s)	Key Algorithmic Features/Approach	Supported Receptor Types	Input Format	Output Format	Strengths	Limitations/Challenges	Integration	Publication/Link
Cell Ranger (10x)	Raw data processing, alignment, assembly, basic clonotyping, GEX integration	Proprietary 10x Genomics workflow, built-in alignment/assembly	TCR, BCR	FASTQ (10x V(D)J + GEX)	.csv, .json, Loupe files	All-in-one solution, user-friendly, tight integration with Loupe browser	Primarily for 10x data, fixed clonotype definition, limited customization	Loupe VDJ Browser, Seurat, Scanpy (via file I/O)	<a href="#">10x Genomics</a>
TRUST4	De novo assembly of TCR/BCR sequences from RNA-seq	de Bruijn graph-based algorithm, supports partial/full CDR3 assembly	TCR, BCR	FASTQ, BAM (sc/bulk RNA-seq)	FASTA, TSV	No V(D)J enrichment needed, works on various RNA-seq types, fast	Accuracy depends on RNA-seq coverage, may miss very low-expression clones	Can be an upstream tool for other software	<a href="#">Li et al., Nat Methods, 2021</a>
TraCeR	Reconstructs paired TCR sequences from scRNA-seq	De novo assembly using Trinity, followed by alignment to IMGT	TCR	FASTQ (scRNA-seq)	TSV, FASTA	Recovers TCRs from non-enriched scRNA-seq, open-source	Computationally intensive (relies on Trinity), limited sensitivity for low-expression TCRs	-	<a href="#">Stubington et al., Nat Methods, 2016</a>
Platypus	Reconstructs TCRs from scRNA-seq, optimized for low expression	Combines k-mer counting and expectation-maximization algorithm	TCR	FASTQ (scRNA-seq)	TSV	Optimized for low TCR expression, increases detection rate	Focuses on TCRs, can be computationally complex	-	<a href="#">Yermanos et al., NAR, 2021</a>
SONG	Fast and accurate V(D)J gene segment alignment	Optimized gapped-alignment algorithm	TCR, BCR	FASTQ, FASTA	SAM-like, custom	High accuracy and speed	Primarily an alignment tool, requires other software for downstream analysis	-	<a href="#">Sheng et al., Bioinformatics, 2022</a>
MiXCR	TCR/BCR alignment, CDR3 extraction, clone quantification	Highly optimized alignment, UMI support, built-in error correction	TCR, BCR	FASTQ, FASTA	TSV, VDJCA	High accuracy, comprehensive, supports various species/protocols	Commercial (free for academia), command-line	VDJtools, Immunarch	<a href="#">Bolotin et al., Nat Methods, 2015</a>
scirpy	Single-cell V(D)J data analysis and integration	Based on Scanpy AnnData, flexible clonotyping, network analysis	TCR, BCR	AnnData, Cell Ranger output, etc.	AnnData, plots	Seamless integration with Scanpy, rich features, flexible, powerful visualization	Requires Python/Scanpy environment, learning curve for non-Python users	Scanpy	<a href="#">Sturm et al., Nat Commun, 2021</a>

Tool Name Dandelion	Primary Functionality Single-cell V(D)J analysis and integration (BCR focus)	Key Algorithmic Features/Approach emphasizes BCR lineage tree construction	Supported Receptor Types (esp. BCR)	Input Format Cell Ranger output	Output Format AnnData, Newick	Strengths Integrates with Scanpy, focuses on BCR lineage analysis	Limitations/Challenges Relatively new, smaller user community than scirpy	Integration Scanpy	Publication/Link Suo et al., Nat Commun, 2022
VDJtools	Post-processing and statistical analysis of repertoire data	Java-based, supports many input formats, rich statistical functions	TCR, BCR	MiXCR, IMGT, Cell Ranger, etc.	TSV, plots	Wide format support, comprehensive statistics, cross-platform	Command-line, non-interactive plots, doesn't process raw data	MiXCR, Immunarch	Shugay et al., PLoS Comput Biol, 2015
Immcantation	Framework for BCR/antibody sequence analysis	R package suite, focuses on BCR SHM, clonal lineage, etc.	BCR (mainly), TCR	FASTA, AIRR Standard	AIRR Standard, R objects	Extremely powerful and mature for BCR analysis, follows AIRR standards	R-based, learning curve for non-R users	R ecosystem, AIRR community tools	immcantation.org
Loupe VDJ Browser	Interactive visualization for 10x Genomics V(D)J data	Desktop application	TCR, BCR	10x .cloupe files	-	User-friendly, highly interactive, easy to link VDJ and GEX	Supports only 10x data format, limited analysis functions	Cell Ranger	10x Genomics

## IV. Navigating Public Single-Cell V(D)J Dataset Resources

The accumulation and sharing of public datasets are crucial for advancing immune repertoire research. They not only provide valuable resources for the development and benchmarking of computational methods but also offer opportunities for bioinformatics analysis and scientific discovery to researchers who cannot generate large-scale data themselves.

### A. The Value of Public Data Repositories

Public datasets serve as a catalyst for accelerating the entire field by promoting meta-analysis, providing benchmarks for tool development, empowering researchers without wet-lab capabilities, and fostering reproducible research.

- Facilitating Meta-analysis and Discovery of Universal Immunological Principles:** By integrating data from different studies, we can test the universality of immunological findings across different populations, disease contexts, or experimental conditions.
- Providing Benchmark Datasets for Tool Development:** New bioinformatics algorithms need to be tested and validated using datasets with known characteristics or thorough annotation.
- Empowering Researchers without Wet-Lab Capabilities:** Computational biologists can use public data for purely in silico research.
- Promoting Reproducibility and Data Sharing:** Adhering to the FAIR principles (Findable, Accessible, Interoperable, Reusable), making data public helps to improve transparency and reproducibility.

### B. Major Data Portals and Consortia for Immune Repertoire Data

Currently, data is primarily stored in general-purpose repositories (like GEO) and specialized immune repertoire databases (like iReceptor, VDJdb), while large research consortia (like HCA) are also major data sources.

#### 1. General-Purpose Repositories:

- NCBI Gene Expression Omnibus (GEO) and Sequence Read Archive (SRA):** The main platforms for sharing all kinds of high-throughput sequencing data, including single-cell V(D)J.
- European Genome-phenome Archive (EGA):** Primarily stores genotype and phenotype data with individual identifiability, requiring authorized access.

#### 2. Specialized Immune Repertoire Databases:

- iReceptor:** A portal for querying AIRR-compliant repertoire data from distributed repositories.
- VDJdb:** A curated database of TCR sequences with known antigen specificity.
- McPAS-TCR:** A catalog of TCR sequences associated with various pathologies and specific antigens.
- TBAdb:** A database of T-cell and B-cell receptors and their corresponding antigens.
- Observed Antibody Space (OAS):** A large-scale collection of published paired and unpaired BCR/antibody sequences.
- Immune Epitope Database (IEDB):** A comprehensive resource containing links to TCR and BCR data recognizing specific epitopes.

#### 3. Data from Large Research Consortia or Specific Projects:

- Human Cell Atlas (HCA):** An international project to create a comprehensive reference map of all human cell types, including a vast amount of single-cell V(D)J data.
- Disease-Specific Research Initiatives:** For example, extensive research on COVID-19 has generated massive public V(D)J datasets.

### C. Considerations When Using Public Datasets

When using public data, it is crucial to be aware of differences in processing pipelines and clonotype definitions across studies and to guard against batch effects to ensure the accuracy and comparability of results.

- Differences in Processing and Clonotype Definitions:** Different studies may use different bioinformatics workflows, affecting cross-study comparisons.
- Importance of Metadata Quality:** High-quality, detailed metadata (species, tissue, disease state, cell type, etc.) are crucial for correct data interpretation.
- Adherence to Data Standards:** Data following standards like AIRR greatly enhance usability and interoperability.
- Batch Effects:** When combining data from different studies, potential batch effects must be identified and properly handled.

The value of public databases is directly related to the standardization of their data and the richness of their metadata, which is essential for enabling large-scale meta-analyses and driving the development of predictive algorithms.

Raw sequence data alone, without contextual information, has limited scientific value. For instance, a researcher looking for TCRs reactive to a specific viral peptide needs more than just the TCR sequence; they need to know the species, tissue of origin, host HLA type, T-cell functional state, and whether the specificity has been experimentally validated. If this metadata is missing or inconsistent across databases, large-scale meta-analysis becomes extremely difficult. The iReceptor Gateway, by promoting the AIRR standards, is attempting to solve this problem.

The availability of public datasets, especially those linking to antigen specificity, has greatly fueled the development of machine learning prediction algorithms.

Machine learning models require large amounts of labeled data for training. Databases that provide TCR sequence-epitope pairings serve as this training data. As more high-quality, experimentally validated data becomes publicly available, the performance of these predictive models will continue to improve, potentially accelerating the identification of therapeutic TCRs.

Public datasets are not only a resource for immunologists but also provide valuable raw material for computational biologists to develop new algorithms.

The unique structure of V(D)J data (e.g., paired chain information, linkage to the transcriptome) presents interesting challenges and opportunities for developing new algorithms. For example, defining cell identity based on both transcriptome and receptor sequence, or tracing B-cell lineages using somatic hypermutation in BCRs, are active areas of computational research. Public datasets provide the raw material for these methodological innovations.

Table 2: Examples of Available Public Single-Cell V(D)J Datasets

*Note: This is an exemplary overview. The availability and specific content of public datasets are constantly updated. Users should consult the relevant databases and publications for the most current and detailed information.*

Dataset ID/Name	Primary Publication/Consortium	Species	Tissue/Cell Type(s)	Disease/Condition	Focus of Study	Single-Cell Tech.	Approx. # of Cells	Key V(D)J-related Findings/Relevance	Data Accession/Link
GSE150728	Wu et al., Cell, 2020	Human	PBMCs	COVID-19 (Severe vs. Mild)	SARS-CoV-2-specific T/B cell response	10x 5' scRNA-seq + VDJ	~150K cells	Revealed highly expanded yet potentially exhausted CD8+ T cells in severe	GEO: <a href="#">GSE150728</a>

Dataset ID/Name	Primary Publication/Consortium	Species	Tissue/Cell Type(s)	Disease/Condition	Focus of Study	Single-cell RNA-seq + VDJ, SMART-seq2	Approx. # of Cells	Key V(D)J-related Findings/Relevance	Data Accession/Link
GSE125970	Azizi et al., Cell, 2018	Human	Tumor, LN, Blood	Breast Cancer	T-cell states and TCR repertoire in TME	TCR-A-seq + VDJ, SMART-seq2	~45K T cells	Provides baseline transcriptome and TCR features of breast cancer T-cell subsets.	<a href="#">GEO: GSE125970</a>
HCA (various)	HCA Consortium	Human	Various healthy tissues	Healthy	Building a cell atlas of healthy human tissues	Various (10x VDJ)	Millions	Provides baseline TCR/BCR repertoire data from various healthy tissues.	<a href="#">HCA Data Portal</a>
GSE110681	Guo et al., Nature, 2018	Human	Colorectal Cancer, Blood	Colorectal Cancer	TCR repertoire & neoantigen response of TILs	SMART-seq2 + TCRseq	~11K T cells	Identified tumor-specific CD8+ T-cell clones and their link to neoantigen reactivity.	<a href="#">GEO: GSE110681</a>
VDJdb	Bagaev et al., NAR, 2020	Human, Mouse	Various	Various	Curated database of TCRs with known antigen specificity	Various	>70K TCR entries	A valuable resource of validated TCR-antigen pairs for research and tool development.	<a href="#">vdjdb.cdr3.net</a>
OAS database	Olsen et al., NAR, 2022	Human, Mouse	Various	Various	Large-scale collection and standardization of BCR/antibody sequences	Various	>1.8B sequences	Massive BCR sequence data supporting antibody engineering and ML models.	<a href="#">opentargets.github.io/oasis/</a>
COVID-19 Cell Atlas	Various	Human	PBMC, BALF, Lung, etc.	COVID-19	Immunopathology of COVID-19, T/B cell response	Various (with VDJ)	Large collection	Integrates many COVID-19 datasets to study repertoire dynamics post-infection.	<a href="#">covid19cellatlas.org</a>
iReceptor Gateway	Corrie et al., Front Immunol, 2018	Various	Various	Various	Federated query of AIRR-compliant repertoire data	N/A	N/A	Facilitates cross-study data retrieval and comparison via AIRR standards.	<a href="#">gateway.ireceptor.org</a>

## V. Connecting Sequence to Insight: Linking Repertoire Features with Cellular States and Biological Context

The true power of single-cell V(D)J sequencing lies in its ability to connect the sequence information of immune receptors with other molecular features of the same cell and with the broader biological context. This integrated analysis is the key step from descriptive repertoire analysis to a mechanistic understanding of immunology.

### A. The Power of Integrated Analysis

**The power of single-cell V(D)J analysis lies in its ability to directly link clonal properties with the phenotype of the same cell, thereby elevating research from sequence cataloging to an investigation of clonal functional roles.** The core strength of modern single-cell V(D)J sequencing methods is the ability to simultaneously capture V(D)J sequences and other cellular features from the same cell. This intrinsic linkage allows researchers to directly correlate clonal properties (e.g., a specific CDR3 sequence, clonal expansion) with cellular phenotypes (e.g., activation status, exhaustion markers, cytokine expression profile, cell type identity). This capability elevates the study of immune repertoires from mere sequence cataloging to an investigation of the functional roles of specific clones in specific biological contexts. For example, knowing that a particular TCR clone is expanded in a tumor sample is interesting, but discovering through integrated transcriptomic data that this expanded TCR clone is predominantly found in cells that highly express PD-1, LAG-3, and TOX (classic markers of T-cell exhaustion) provides a much richer and more actionable insight. This might imply that although this clone has expanded, its function may already be compromised, thus suggesting a potential strategy of applying immune checkpoint inhibitors (like anti-PD-1 antibodies) to "reawaken" these T-cells.

### B. Linking Clonotypes with Transcriptomic Features

By integrating transcriptomic data, one can identify gene expression programs and molecular pathways associated with specific clonal behaviors, thereby revealing their functional state. By integrating V(D)J data with scRNA-seq data from the same cells, one can:

- Identify gene expression programs associated with expanded clones:** For example, after cancer immunotherapy, one can analyze whether significantly expanded T-cell clones exhibit transcriptional signatures of effector T-cells (e.g., high expression of *GZMB*, *IFNG*), memory T-cells (e.g., high expression of *IL7R*, *CCR7*), or exhausted T-cells (e.g., high expression of *PDCD1*, *CTLA4*, *LAG3*, *TOX*).
- Perform differential gene expression analysis:** By grouping cells according to their clonotype (or clonotype size) and then comparing the gene expression profiles between groups, one can discover molecular pathways and regulatory networks associated with specific clonal behaviors.

Patterns of association between V-gene usage and cell activation states discovered from integrated data can drive new hypotheses about the rules of immune recognition and response.

Patterns emerging from integrated single-cell data, such as a specific V-gene usage consistently appearing with an activated T-cell phenotype (e.g., high Ki67 expression) after vaccination, can drive new hypotheses. If T-cells using Vβ5.1 and expressing Ki67 consistently expand across multiple individuals vaccinated with vaccine X, it might suggest that the Vβ5.1 gene itself is particularly well-suited for recognizing a dominant epitope in vaccine X.

### C. Linking Clonotypes with Cell Surface Protein Expression (CITE-seq)

CITE-seq technology provides a more direct, quantitative measurement of cell surface proteins, enabling a more refined phenotypic analysis and functional state confirmation of immune cells. CITE-seq (Cellular Indexing of Transcriptomes and Epitopes by Sequencing) technology, using oligonucleotide-labeled antibodies, makes it possible to simultaneously measure RNA expression and cell surface protein levels at the single-cell level.

- Refined immunophenotyping:** Directly linking specific TCR/BCR sequences with detailed immune phenotypes (e.g., distinguishing naive/memory T-cell subsets using CD45RO, CCR7; assessing T-cell exhaustion with PD-1, TIM-3, LAG-3).
- Functional state confirmation:** Surface protein expression can provide corroborating evidence or a more precise definition for functional states inferred from transcriptomic data.

### D. Tracking Clonal Dynamics Over Time or Conditions

Through longitudinal and comparative studies, clonal dynamics can be tracked over time or across conditions, leading to the discovery of immunological markers associated with clinical outcomes.

- Longitudinal studies:** Periodically collecting samples during infection, after vaccination, or during cancer treatment for single-cell V(D)J and multi-omics analysis allows for observing how specific clonotypes expand, contract, or how their phenotypes and functional states evolve over time.
- Comparative studies:** Comparing immune repertoire features and associated cellular states between different patient cohorts (e.g., responders vs. non-responders to a treatment) can help discover immunological markers associated with clinical outcomes.

Linking trackable, specific clonotypes to clinical outcomes (like treatment response, disease progression, or vaccine protection) is the most translationally potent area of single-cell V(D)J analysis. If a set of TCR clonotypes consistently persists in cancer patients who respond to immunotherapy but is absent in non-responders, these clonotypes themselves could become biomarkers for predicting treatment response.

### E. Challenges in Establishing Causality

Although integrated analysis can reveal strong correlations, establishing causality typically requires further functional experiments, a critical step from observation to mechanistic understanding.

- Correlation does not equal causation:** For example, a T-cell clone that is expanded and shows an exhausted phenotype in a tumor might have expanded and then become exhausted due to chronic antigen exposure, or the clone itself might have a predisposition to exhaustion.
- Necessity of functional validation:** For TCRs or BCRs identified as biologically significant in integrated analyses, their properties usually need to be validated through functional experiments, such as cloning the



sequence into a reporter cell line to test its reactivity to antigens.

## VI. Emerging Frontiers and Future Outlook

Single-cell V(D)J sequencing and its analysis methodologies are in a state of rapid development. New technologies, new algorithms, and new biological questions are constantly emerging and shaping the future of the field.

### A. Technological Advances

#### 1. Spatial Context Integration:

**Integrating single-cell analysis with spatial transcriptomics is a current technological frontier, enabling the in-situ resolution of immune clone locations in tissues, which is crucial for understanding cell-cell interactions.** This will allow researchers to understand the precise spatial localization of clonotypes with specific TCRs/BCRs within the tissue microenvironment. For instance, a cytotoxic T-cell clone may be abundant in a tumor tissue, but its anti-tumor effect will be greatly diminished if it is confined to the stromal region and cannot effectively infiltrate the tumor bed.

#### 2. Higher-Throughput Multi-Omics:

**The future trend is to integrate more dimensions of omics information at the single-cell level to provide a more comprehensive portrait of individual immune cell states.** This includes combining V(D)J sequences and transcriptomes with ATAC-seq (for chromatin accessibility), metabolomics, and more.

#### 3. Application of Long-Read Sequencing in V(D)J Analysis:

**Long-read sequencing holds the promise of more directly obtaining complete V(D)J sequences, but its throughput and cost in single-cell applications remain challenges to be addressed.** While current mainstream single-cell V(D)J sequencing relies on short-read technology, long-read technologies (like PacBio or Oxford Nanopore) can theoretically obtain full-length, un-spliced TCR/BCR transcripts more directly, but their throughput, cost, and error rate for single-cell applications need continuous improvement.

### B. Computational Challenges and Innovations

#### 1. TCR/BCR-Antigen Specificity Prediction:

**Predicting TCR/BCR antigen specificity using AI is the "Holy Grail" of computational immunology; its success would revolutionize the development of personalized immunotherapies and diagnostics.**

Using machine learning and AI models to predict the antigenic epitopes recognized by a TCR or BCR based on its sequence and/or structural information is one of the "Holy Grails" of the field. This requires large, high-quality, experimentally validated receptor-antigen interaction data for training. If successful, it would allow us to rapidly identify disease-relevant receptors purely from a patient's repertoire sequence data, greatly accelerating the development of diagnostics and personalized treatment strategies.

#### 2. Dynamic Modeling of Repertoire Evolution:

**Developing computational models that can simulate and predict the dynamic evolution of the immune repertoire is an important direction for understanding the complex dynamics of the immune system.** This requires developing complex models that can integrate multiple factors and make forward-looking predictions.

#### 3. Application of Advanced Machine Learning and AI:

**Beyond specificity prediction, AI/ML can be used to identify complex immune signatures relevant to clinical outcomes and to integrate heterogeneous multi-omics data.** This includes identifying subtle but critical "immune signatures" from high-dimensional repertoire data that correlate with disease state or treatment response.

#### 4. Improved Clonotype Definition and Lineage Tracing:

**Developing more refined methods for clonotype definition and lineage tracing, especially for complex BCR evolution, is a direction of continuous innovation in computational methodology.** This is particularly true for BCRs, where more accurate and efficient tools are needed to trace the evolutionary lineages of B-cells during affinity maturation.

### C. Expanding Biological Understanding

Future biological research will focus on a deeper understanding of cross-reactivity, the interplay between innate and adaptive immunity, and ultimately, the realization of personalized immunotherapy and vaccine design.

- **Deepening the understanding of TCR/BCR cross-reactivity:** Unveiling the molecular mechanisms and sequence/structural features that determine cross-recognition ability.
- **Elucidating the interplay between innate and adaptive immunity:** Investigating the mutual regulation between different immune cell types through integrated single-cell multi-omics.
- **Personalized immunotherapy and vaccine design:** Designing more precise and effective personalized therapies and vaccines based on individual repertoire characteristics.

### D. Ethical Considerations

As immune repertoire data becomes increasingly powerful, ethical issues such as data privacy, secure sharing, and the potential for discrimination urgently require the establishment of responsible guidelines.

- **Data Privacy and Identifiability:** An individual's immune repertoire is highly unique and could potentially be used for re-identification.
- **Data Security and Sharing Protocols:** Strong anonymization techniques, secure data storage, and clear ethical and legal frameworks are needed.
- **Potential for Discrimination:** There is a need to guard against the misuse of information about immune repertoire features in areas such as insurance or employment.

## VII. Conclusion

**Single-cell V(D)J sequencing and its multi-omics integration strategies have become revolutionary tools for insight into adaptive immunity, with their development relying on the synergistic progress of technology, computation, and data resources.** Single-cell V(D)J sequencing and its related multi-omics integration strategies have provided revolutionary tools for us to gain insight into the complexity of the adaptive immune response with unprecedented resolution and depth. By simultaneously resolving the antigen receptor sequence, transcriptome profile, surface protein expression, and even spatial location of individual immune cells, researchers can more accurately identify functional immune clones, track their dynamic changes, and tightly link these molecular features to the biological behavior of cells and clinical phenotypes.

**Extracting meaningful biological insights from raw data is highly dependent on complex, sophisticated computational analysis methods and robust bioinformatics tools.** This report has reviewed the cutting-edge methods in key computational steps, from data pre-processing, V(D)J sequence annotation, and clonotype definition to the quantification of repertoire features, visualization, and multi-modal data integration. It has also emphasized the importance of standardized data formats (like the AIRR standards) and high-quality public datasets (such as iReceptor, VDJdb, GEO) for promoting the field's development, enabling data sharing, and facilitating meta-analysis.

**The field is in a virtuous cycle of technological innovation, computational method advancement, and data resource growth; this synergistic interaction is the core driver of progress in the field.** The progress described throughout this report stems not from a single factor, but from the synergistic interaction between these elements. More advanced sequencing technologies generate more and more complex data, which in turn requires more powerful computational tools for processing and analysis; new analytical methods reveal novel biological patterns, thus inspiring new experimental designs and data generation strategies; and the sharing of public data accelerates this virtuous cycle by allowing multiple researchers to participate in tool development and scientific discovery.

**The ultimate goal of the field is to translate vast amounts of data into actionable immunological knowledge for clinical application, which requires sustained, transdisciplinary collaboration.** Achieving this grand goal requires continuous, transdisciplinary collaboration among basic scientists, computational biologists, clinicians, and ethicists. Single-cell V(D)J sequencing is undoubtedly a powerful engine driving this translational process, and its future development and application prospects are vast, promising to bring profound changes to multiple fields, including infectious diseases, oncology, autoimmune diseases, vaccine development, and personalized medicine.

tip: click on the pencil icon on the left to clear the editor

## GitHub flavoured styling by default

We now use GitHub flavoured styling by default.