

# Package ‘immApex’

July 16, 2025

**Title** Tools for Adaptive Immune Receptor Sequence-Based Machine and Deep Learning

**Version** 1.3.0

**Description** A set of tools to build tensorflow/keras3-based models in R from amino acid and nucleotide sequences focusing on adaptive immune receptors. The package includes pre-processing of sequences, unifying gene nomenclature usage, encoding sequences, and combining models. This package will serve as the basis of future immune receptor sequence functions/packages/models compatible with the scRepertoire ecosystem.

**License** MIT + file LICENSE

**Encoding** UTF-8

**RoxygenNote** 7.3.2

**biocViews** Software, ImmunoOncology, SingleCell, Classification, Annotation, Sequencing, MotifAnnotation

**Depends** R (>= 4.3.0)

**Imports** hash, httr, igraph, keras3, magrittr, matrixStats, methods, Rcpp (>= 0.12.11), reticulate, rvest, SingleCellExperiment, stats, stringi, stringr, tensorflow, utils

**Suggests** BiocStyle, ggraph, ggplot2, knitr, graph, markdown, rmarkdown, scRepertoire, spelling, testthat, tidygraph, viridis

**LinkingTo** Rcpp

**VignetteBuilder** knitr

**Language** en-US

**URL** <https://github.com/BorchLab/immApex/>

**BugReports** <https://github.com/BorchLab/immApex/issues>

**git\_url** <https://git.bioconductor.org/packages/immApex>

**git\_branch** devel

**git\_last\_commit** f465914

**git\_last\_commit\_date** 2025-05-22

**Repository** Bioconductor 3.22

**Date/Publication** 2025-07-16

**Author** Nick Borcharding [aut, cre]

**Maintainer** Nick Borcharding <ncborch@gmail.com>

## Contents

immApex-package	2
adjacencyMatrix	3
buildNetwork	3
formatGenes	5
generateSequences	6
geometricEncoder	7
getIMGT	8
getIR	9
immapex_AA.data	9
immapex_blosum.pam.matrices	10
immapex_example.data	11
immapex_gene.list	11
inferCDR	12
mutateSequences	13
onehotEncoder	14
positionalEncoder	15
probabilityMatrix	15
propertyEncoder	16
sequenceDecoder	17
tokenizeSequences	19
variationalSequences	20
<b>Index</b>	<b>22</b>

---

immApex-package	<i>immApex: Tools for Adaptive Immune Receptor Sequence-Based Machine and Deep Learning</i>
-----------------	---

---

## Description

A set of tools to build tensorflow/keras3-based models in R from amino acid and nucleotide sequences focusing on adaptive immune receptors. The package includes pre-processing of sequences, unifying gene nomenclature usage, encoding sequences, and combining models. This package will serve as the basis of future immune receptor sequence functions/packages/models compatible with the scRepertoire ecosystem.

## Author(s)

**Maintainer:** Nick Borchering <ncborch@gmail.com>

## See Also

Useful links:

- <https://github.com/BorchLab/immApex/>
- Report bugs at <https://github.com/BorchLab/immApex/issues>

---

adjacencyMatrix*Adjacency matrix from amino acid or nucleotide sequences*

---

**Description**

Calculate frequency of adjacency between residues along a set of biological sequences.

**Usage**

```
adjacencyMatrix(  
  input.sequences = NULL,  
  normalize = TRUE,  
  sequence.dictionary = amino.acids  
)
```

**Arguments**

input.sequences	The amino acid or nucleotide sequences to use
normalize	Return the values as a function of total number of residues ( <b>TRUE</b> ) or frequencies ( <b>FALSE</b> )
sequence.dictionary	The letters to use in sequence generation (default are all amino acids)

**Value**

Adjacency matrix based on input.sequences.

**Examples**

```
new.sequences <- generateSequences(prefix.motif = "CAS",  
                                   suffix.motif = "YF",  
                                   number.of.sequences = 100,  
                                   min.length = 8,  
                                   max.length = 16)  
  
adj.matrix <- adjacencyMatrix(new.sequences,  
                              normalize = TRUE)
```

---

buildNetwork*Build Edit Distance Network Using Symmetric Deletion Lookup*

---

**Description**

Constructs a weighted similarity network from biological sequences using a symmetric deletion lookup strategy combined with a banded edit-distance computation. The returned igraph object contains vertices representing the input sequences and edges representing pairs of sequences whose edit distance is less than or equal to the specified threshold. The edge attribute weight stores the computed edit distance.

## Usage

```
buildNetwork(
  input.data,
  sequence.column = "sequence",
  threshold = 2,
  filter.v = FALSE,
  filter.j = FALSE,
  technology = NULL,
  simplify.format = TRUE,
  simplify.families = TRUE
)
```

## Arguments

<code>input.data</code>	A character vector of AIR sequences, or a data frame containing sequence data.
<code>sequence.column</code>	A character string specifying the name of the column in <code>input.data</code> that contains the sequences. Default is "sequence". This parameter is ignored when <code>input.data</code> is a character vector.
<code>threshold</code>	An integer specifying the maximum allowed edit distance. Only pairs of sequences with an edit distance less than or equal to this value will be connected. Default is 2.
<code>filter.v</code>	Logical indicating whether to filter candidate pairs to only those that have matching <code>v.gene</code> family annotations. Default is FALSE. When TRUE, the input data frame must contain a column with V gene annotations, either named <code>v.gene</code> or determined by <code>.get.genes.updated</code> .
<code>filter.j</code>	Logical indicating whether to filter candidate pairs to only those that have matching <code>j.gene</code> family annotations. Default is FALSE. When TRUE, the input data frame must contain a column with J gene annotations, either named <code>j.gene</code> or determined by <code>.get.genes.updated</code> .
<code>technology</code>	The sequencing technology employed - 'TenX', 'Adaptive', or 'AIRR'.
<code>simplify.format</code>	If applicable, remove the allelic designation ( <b>TRUE</b> ) or retain all information ( <b>FALSE</b> )
<code>simplify.families</code>	If applicable, remove the hyphenated designation ( <b>TRUE</b> ) or retain all information ( <b>FALSE</b> )

## Details

This function supports both a character vector of sequences and a data frame. When provided a data frame, the user can specify the column containing sequences using the `sequence.column` parameter. Additionally, candidate pairs can be filtered by requiring matching `v.gene` and/or `j.gene` annotations (see `filter.v` and `filter.j`). If filtering is enabled, the corresponding gene annotation columns are required.

The function first calls a C++ routine (via Rcpp) to perform a symmetric deletion lookup, generating candidate pairs of sequences that might be within the specified edit distance. It then uses a banded dynamic programming algorithm (also implemented in C++) to compute the exact edit distance for each candidate pair. When using a data frame input, the candidate pairs can be further filtered by requiring that sequences have matching `v.gene` and/or `j.gene` values. Note that gene filtering is only applied if the corresponding filtering flag is set to TRUE.

**Value**

An igraph object representing the AIR similarity network. Vertices contain the original sequences (and gene annotations, if available), and each edge has a weight attribute corresponding to the computed edit distance. If no edges meet the threshold, an igraph object with only vertices is returned.

**Examples**

```
# Using a character vector of sequences:
sequences <- c("CASSLGTDQYF", "CASSPGTDQYF", "CASSLGNDQYF", "CASRLGNDQYF")
g <- buildNetwork(sequences, threshold = 2)
plot(g)

# Using a data frame with a custom sequence column:
df <- data.frame(
  mySeqs = c("CASSLGTDQYF", "CASSPGTDQYF", "CASSLGNDQYF", "CASRLGNDQYF"),
  v = c("TRBV20", "TRBV20", "TRBV12", "TRBV20"),
  j = c("TRBJ2-7", "TRBJ2-7", "TRBJ2-1", "TRBJ2-7")
)
g_df <- buildNetwork(df,
  threshold = 2,
  filter.v = TRUE,
  filter.j = TRUE,
  sequence.column = "mySeqs")
plot(g_df)
```

formatGenes

*Ensure clean gene nomenclature using IMGT annotations***Description**

This function will format the genes into a clean nomenclature using the IMGT conventions.

**Usage**

```
formatGenes(
  input.data,
  region = "v",
  technology = NULL,
  species = "human",
  simplify.format = TRUE
)
```

**Arguments**

input.data	Data frame of sequencing data or scRepertoire outputs
region	Sequence gene loci to access - "v", "d", "j", or "c" or a combination using c("v", "d", "j")
technology	The sequencing technology employed - 'TenX', 'Adaptive', or 'AIRR'

**species** One or two word designation of species. Currently supporting: "human", "mouse", "rat", "rabbit", "rhesus monkey", "sheep", "pig", "platypus", "alpaca", "dog", "chicken", and "ferret"

**simplify.format** If applicable, remove the allelic designation (**TRUE**) or retain all information (**FALSE**)

### Value

A data frame with the new columns of formatted genes added.

### Examples

```
data(immapex_example.data)
formatGenes(immapex_example.data[["TenX"]],
            region = "v",
            technology = "TenX")
```

---

generateSequences	<i>Randomly Generate Amino Acid Sequences</i>
-------------------	---

---

### Description

Use this to make synthetic amino acid sequences for purposes of testing code, training models, or noise.

### Usage

```
generateSequences(
  prefix.motif = NULL,
  suffix.motif = NULL,
  number.of.sequences = 100,
  min.length = 1,
  max.length = 10,
  sequence.dictionary = amino.acids
)
```

### Arguments

**prefix.motif** Add defined amino acid/nucleotide sequence to the start of the generated sequences.

**suffix.motif** Add defined amino acid/nucleotide sequence to the end of the generated sequences

**number.of.sequences** Number of sequences to generate

**min.length** Minimum length of the final sequence. The min.length may be adjusted if incongruent with prefix.motif/suffix.motif lengths

**max.length** Maximum length of the final sequence

**sequence.dictionary** The letters to use in sequence generation (default are all amino acids)

**Value**

A vector of generated sequences

**Examples**

```
generateSequences(prefix.motif = "CAS",
                  suffix.motif = "YF",
                  number.of.sequences = 100,
                  min.length = 8,
                  max.length = 16)
```

---

geometricEncoder

*Geometric Encoder from Amino Acid Strings*


---

**Description**

Use this to transform amino acid sequences into a geometric encoding of the sequence.

**Usage**

```
geometricEncoder(
  input.sequences,
  method.to.use = "BLOSUM62",
  theta = pi/3,
  verbose = TRUE
)
```

**Arguments**

input.sequences	The set of amino acid sequences
method.to.use	The method or approach for the conversion: <ul style="list-style-type: none"> <li>• BLOcks SUBstitution Matrices: BLOSUM45, BLOSUM50, BLOSUM62, BLOSUM80, BLOSUM100</li> <li>• Point Accepted Mutation Matrices: PAM30, PAM40, PAM70, PAM120, PAM250</li> </ul>
theta	The angle for geometric transformation
verbose	Print messages corresponding to the processing step

**Value**

Geometric encoded amino acid sequences in a matrix

## Examples

```
new.sequences <- generateSequences(prefix.motif = "CAS",
                                   suffix.motif = "YF",
                                   number.of.sequences = 100,
                                   min.length = 8,
                                   max.length = 16)

sequence.matrix <- geometricEncoder(new.sequences,
                                    method.to.use = "BLOSUM62",
                                    theta = pi/3)
```

---

getIMGT

---

*Get IMGT Sequences for Specific Loci*


---

## Description

Use this to access the ImMunoGeneTics (IMGT) sequences for a specific species and gene loci. More information on IMGT can be found at [imgt.org](http://imgt.org).

## Usage

```
getIMGT(
  species = "human",
  chain = "TRB",
  sequence.type = "aa",
  frame = "inframe",
  region = "v",
  max.retries = 3,
  verbose = TRUE
)
```

## Arguments

<code>species</code>	One or two-word common designation of species.
<code>chain</code>	Sequence chain to access, e.g., <b>TRB</b> or <b>IGH</b> .
<code>sequence.type</code>	Type of sequence - <b>aa</b> (amino acid) or <b>nt</b> (nucleotide).
<code>frame</code>	Designation for <b>all</b> , <b>inframe</b> , or <b>inframe+gap</b> .
<code>region</code>	Gene loci to access.
<code>max.retries</code>	Number of attempts to fetch data in case of failure.
<code>verbose</code>	Print messages corresponding to the processing step.

## Value

A list of allele sequences.



**Examples**

```
TRBV_aa <- getIMGT(species = "human",
                   chain = "TRB",
                   frame = "inframe",
                   region = "v",
                   sequence.type = "aa",
                   max.retries = 3)
```

getIR

*Extract Immune Receptor Sequences***Description**

Use this to extract immune receptor sequences from a Single-Cell Object or the output of [combineTCR](#) and [combineBCR](#).

**Usage**

```
getIR(input.data, chains, sequence.type = "aa")
```

**Arguments**

input.data	Single-cell object or the output of <a href="#">combineTCR</a> and <a href="#">combineBCR</a> from scRepertoire
chains	Immune Receptor chain to use - <b>TRA</b> , <b>TRB</b> , <b>IGH</b> , or <b>IGL</b>
sequence.type	Extract amino acid ( <b>aa</b> ) or nucleotide ( <b>nt</b> ) sequences

**Value**

A data frame of nucleotide or amino acid sequences

immapex\_AA.data

*A list of amino acid properties***Description**

A list of amino acid properties that are used for [propertyEncoder](#) function.

This includes:

- atchleyFactors
- crucianiProperties
- FASGAI
- kideraFactors
- MSWHIM
- ProtFP
- stScales
- tScales
- VHSE
- zScales

**Usage**

```
data("immapex_AA.data")
```

**Value**

List of 10 amino acid properties for 20 amino acids

---

```
immapex_blosum.pam.matrices
```

*List of amino acid substitution matrices*

---

**Description**

A list of amino acid substitution matrices, using the Point Accepted Matrix (PAM) and BLOck SUBstitution Matrix (BLOSUM) approaches. A discussion and comparison of these matrices are available at [PMID: 21356840](#).

- BLOSUM45
- BLOSUM50
- BLOSUM62
- BLOSUM80
- BLOSUM100
- PAM30
- PAM40
- PAM70
- PAM120
- PAM250

**Usage**

```
data("immapex_blosum.pam.matrices")
```

**Value**

List of 10 substitution matrices

---

immapex\_example.data     *Example contig data for Apex*


---

### Description

Contains a collection of bulk or paired TCR sequences in the respective formats in the form of a list from the following sources:

- TenX: 10k\_Human\_DTC\_Melanoma\_5p\_nextgem\_Multiplex from [10x Website](#).
- AIRR: Human\_colon\_16S8157851 from [PMID: 37055623](#).
- Adaptive: Adaptive\_2283\_D0 from [PMID: 36220826](#).

More information on the data formats are available: [AIRR](#), [Adaptive](#), and [TenX](#).

### Usage

```
data("immapex_example.data")
```

### Value

List of 3 example data sets for 10x, AIRR and Adaptive contigs.

---

immapex\_gene.list     *A list of IMGT gene names by genes, loci, and species*


---

### Description

A list of regularized gene nomenclature to use for converting for data for uniformity. Data is organized by gene region, loci and species. Not all species are represented in the data and pseudogenes have not been removed.

### Usage

```
data("immapex_gene.list")
```

### Value

List of gene nomenclature by region, loci, and species.

inferCDR

*Infer portions of the CDR Loop from Vgene data***Description**

Use this isolate sequences from the CDR loop using the V gene annotation. When there are multiple V gene matches for a single gene, the first allelic sequence is used.

**Usage**

```
inferCDR(
  input.data,
  reference = NULL,
  chain = "TRB",
  technology = NULL,
  sequence.type = "aa",
  sequences = c("CDR1", "CDR2")
)
```

**Arguments**

input.data	Data frame output of <a href="#">formatGenes</a>
reference	IMGT reference sequences from <a href="#">getIMGT</a>
chain	Sequence chain to access, like <b>TRB</b> or <b>IGH</b>
technology	The sequencing technology employed - <b>TenX</b> , <b>Adaptive</b> , <b>AIRR</b> , or <b>Omniscope</b>
sequence.type	Type of sequence - <b>aa</b> for amino acid or <b>nt</b> for nucleotide
sequences	The specific regions of the CDR loop to get from the data, such as <b>CDR1</b> .

**Value**

A data frame with the new columns of CDR sequences added.

**Examples**

```
# Getting the Sequence Reference
data(immapex_example.data)
TRBV_aa <- getIMGT(species = "human",
  chain = "TRB",
  frame = "inframe",
  region = "v",
  sequence.type = "aa")

# Ensuring sequences are formatted to IMGT
TenX_formatted <- formatGenes(immapex_example.data[["TenX"]],
  region = "v",
  technology = "TenX")

# Inferring CDR loop elements
TenX_formatted <- inferCDR(TenX_formatted,
  chain = "TRB",
  reference = TRBV_aa,
```

```
technology = "TenX",
sequence.type = "aa",
sequences = c("CDR1", "CDR2"))
```

mutateSequences	<i>Randomly Mutate Sequences of Amino Acids</i>
-----------------	---

### Description

Use this to mutate or mask sequences for purposes of testing code, training models, or noise.

## Usage

```
mutateSequences(  
    input.sequences,  
    n.sequences = 1,  
    mutation.rate = 0.01,  
    position.start = NULL,  
    position.end = NULL,  
    sequence.dictionary = amino.acids  
)
```

## Arguments

input.sequences	The amino acid or nucleotide sequences to use
n.sequences	The number of mutated sequences to return
mutation.rate	The rate of mutations to introduce into sequences
position.start	The starting position to mutate along the sequence <b>Default</b> = NULL will start the random mutations at position 1
position.end	The ending position to mutate along the sequence <b>Default</b> = NULL will end the random mutations at the last position
sequence.dictionary	The letters to use in sequence mutation (default are all amino acids)

## Value

A vector of mutated sequences

## Examples

```
sequences <- generateSequences(prefix.motif = "CAS",
                               suffix.motif = "YF",
                               number.of.sequences = 100,
                               min.length = 8,
                               max.length = 16)

mutated_sequences <- mutateSequences(sequences,
                                     n.sequence = 1,
                                     position.start = 3,
                                     position.end = 8)
```



---

positionalEncoder	<i>Adding Position-Specific Information to Sequences</i>
-------------------	--

---

**Description**

Use this calculate positional encoding for recurrent neural networks using sin/cos and position information.

**Usage**

```
positionalEncoder(number.of.sequences, latent.dims = NULL)
```

**Arguments**

number.of.sequences	The number of sequences to generate position information
latent.dims	The number of latent dimensions.

**Value**

A matrix of values

**Examples**

```
position.info <- positionalEncoder(number.of.sequences = 1000,
                                   latent.dims = 64)
```

---

probabilityMatrix	<i>Position Probability Matrix for Amino Acid or Nucleotide Sequences</i>
-------------------	---

---

**Description**

Use this to generate a position-probability or weight matrix for a set of given sequences.

**Usage**

```
probabilityMatrix(
  input.sequences,
  max.length = NULL,
  convert.PWM = FALSE,
  background.frequencies = NULL,
  sequence.dictionary = amino.acids,
  padding.symbol = ".",
  verbose = TRUE
)
```

**Arguments**

input.sequences	The amino acid or nucleotide sequences to use
max.length	Additional length to pad, NULL will pad sequences to the max length of input.sequences
convert.PWM	Convert the matrix into a positional weight matrix using log likelihood
background.frequencies	Provide amino acid or nucleotide frequencies for the positional weight matrix. If NULL, assumes uniform likelihood.
sequence.dictionary	The letters to use in sequence generation (default are all amino acids)
padding.symbol	Symbol to use for padding at the end of sequences
verbose	Print messages corresponding to the processing step

**Value**

A matrix with position specific probabilities or weights

**Examples**

```
new.sequences <- generateSequences(prefix.motif = "CAS",
                                   suffix.motif = "YF",
                                   number.of.sequences = 100,
                                   min.length = 8,
                                   max.length = 16)

PPM.matrix <- probabilityMatrix(new.sequences)
```

---

propertyEncoder

*Encoder from Amino Acid String by Properties*


---

**Description**

Use this to transform amino acid sequences a a matrix by amino acid properties derived from dimensional reduction strategies

**Usage**

```
propertyEncoder(
  input.sequences,
  max.length = NULL,
  method.to.use = NULL,
  convert.to.matrix = TRUE,
  summary.function = NULL,
  padding.symbol = ".",
  verbose = TRUE
)
```



**Arguments**

<code>input.sequences</code>	The amino acid sequences to use
<code>max.length</code>	Additional length to pad, NULL will pad sequences to the max length of <code>input.sequences</code>
<code>method.to.use</code>	The method or approach to use for the conversion: <ul style="list-style-type: none"> <li>• Individual sets: <code>atchleyFactors</code>, <code>crucianiProperties</code>, <code>FASGAI</code>, <code>kideraFactors</code>, <code>MSWHIM</code>, <code>ProtFP</code>, <code>stScales</code>, <code>tScales</code>, <code>VHSE</code>, <code>zScales</code></li> <li>• Multiple Sets: <code>c("atchleyFactors", "VHSE")</code></li> </ul>
<code>convert.to.matrix</code>	Return a matrix (TRUE) or a 3D array (FALSE)
<code>summary.function</code>	Return a matrix that summarize the amino acid method/property Available summaries include: "median", "mean", "sum", variance ("vars"), or Median Absolute Deviation ("mads")
<code>padding.symbol</code>	Symbol to use for padding at the end of sequences
<code>verbose</code>	Print messages corresponding to the processing step

**Value**

Converted amino acid sequences by property in a matrix or 3D array

**Examples**

```
new.sequences <- generateSequences(prefix.motif = "CAS",
                                   suffix.motif = "YF",
                                   number.of.sequences = 100,
                                   min.length = 8,
                                   max.length = 16)

if(reticulate::py_module_available("numpy")) {
  sequence.matrix <- propertyEncoder(new.sequences,
                                     method.to.use = "VHSE",
                                     convert.to.matrix = TRUE)
}
```

---

sequenceDecoder

---

*One Hot Decoder from One Hot Encoded Matrix or 3D Array*


---

**Description**

Use this to transform one hot encoded sequences back into amino acid or nucleotide sequences.





---

variationalSequences    *Generate Similar Sequences using Variational Autoencoder*

---

## Description

Use this to simulate sequences using a variational autoencoder (VAE) and perturbation of the probability distributions.

## Usage

```
variationalSequences(
  input.sequences,
  encoder.function = "onehotEncoder",
  aa.method.to.use = NULL,
  number.of.sequences = 100,
  encoder.hidden.dim = c(128, 64),
  decoder.hidden.dim = NULL,
  latent.dim = 16,
  batch.size = 16,
  epochs = 50,
  learning.rate = 0.001,
  epsilon.std = 1,
  call.threshold = 0.2,
  activation.function = "relu",
  optimizer = "adam",
  disable.eager.execution = FALSE,
  sequence.dictionary = amino.acids,
  verbose = TRUE
)
```

## Arguments

`input.sequences`  
The amino acid or nucleotide sequences to use

`encoder.function`  
The method to prepare the sequencing information - "onehotEncoder" or "propertyEncoder"

`aa.method.to.use`  
The method or approach to use for the conversion:

- Individual sets: `atchleyFactors`, `crucianiProperties`, `FASGAI`, `kideraFactors`, `MSWHIM`, `ProtFP`, `stScales`, `tScales`, `VHSE`, `zScales`
- Multiple Sets: `c("atchleyFactors", "VHSE")`

`number.of.sequences`  
Number of sequences to generate

`encoder.hidden.dim`  
A vector of the neurons to use in the hidden layers for the encoder portion of the model

`decoder.hidden.dim`  
A vector of the neurons to use in the hidden layers for the decoder portion of the model. If `NULL` assumes symmetric autoencoder



# Index

- \* **internal**
  - immApex-package, [2](#)
- adjacencyMatrix, [3](#)
- buildNetwork, [3](#)
- combineBCR, [9](#)
- combineTCR, [9](#)
- formatGenes, [5](#), [12](#)
- generateSequences, [6](#)
- geometricEncoder, [7](#)
- getIMGT, [8](#), [12](#)
- getIR, [9](#)
- immApex (immApex-package), [2](#)
- immApex-package, [2](#)
- immapex\_AA.data, [9](#)
- immapex\_blosum.pam.matrices, [10](#)
- immapex\_example.data, [11](#)
- immapex\_gene.list, [11](#)
- inferCDR, [12](#)
- mutateSequences, [13](#)
- onehotEncoder, [14](#)
- positionalEncoder, [15](#)
- probabilityMatrix, [15](#)
- propertyEncoder, [9](#), [16](#)
- sequenceDecoder, [17](#)
- tokenizeSequences, [19](#)
- variationalSequences, [20](#)