

R4RNA: A R package for RNA visualization and analysis

Daniel Lai

October 29, 2024

Contents

1 R4RNA	1
1.1 Reading Input	1
1.2 Basic Arc Diagram	2
1.3 Multiple Structures	2
1.4 Filtering Helices	3
1.5 Colouring Structures	3
1.6 Overlapping Multiple Structures	4
1.7 Visualizing Multiple Sequence Alignments	5
1.8 Multiple Sequence Alignments with Annotated Arcs	6
1.9 Additional Colouring Methods	6
1.9.1 Colour By Covariation (with alignment as blocks)	6
1.9.2 Colour By Conservation (with custom alignment colours)	7
1.9.3 Colour By Percentage Canonical Basepairs (with custom arc colours)	7
1.9.4 Colour Pseudoknots (with CLUSTALX-style alignment)	8
2 Session Information	8

1 R4RNA

The R4RNA package aims to be a general framework for the analysis of RNA secondary structure and comparative analysis in R, the language so chosen due to its native support for publication-quality graphics, and portability across all major operating systems, and interactive power with large datasets.

To demonstrate the ease of creating complex arc diagrams, a short example is as follows.

1.1 Reading Input

Currently, supported input formats include dot-bracket, connect, bpseq, and a custom “helix” format. Below, we read in a structure predicted by TRANSAT, the known structure obtained from the RFAM database.

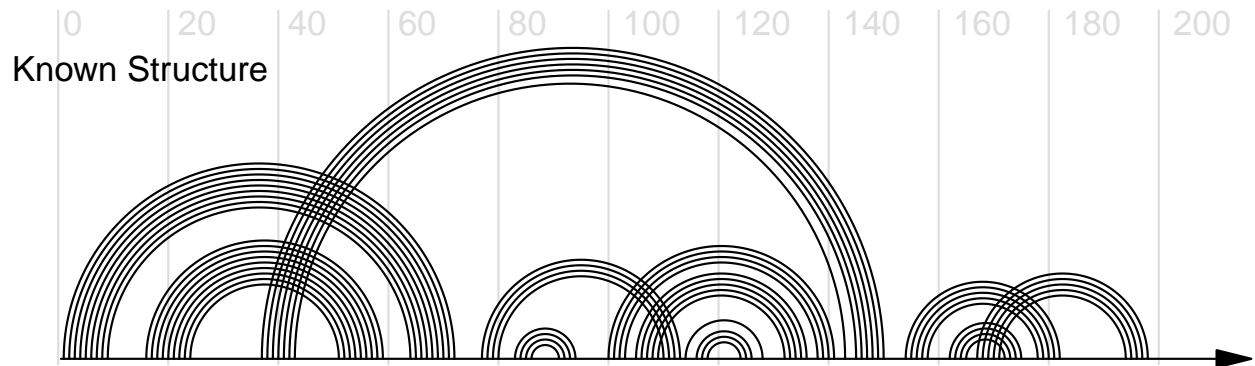
```
> library(R4RNA)
> message("TRANSAT prediction in helix format")
> transat_file <- system.file("extdata", "helix.txt", package = "R4RNA")
> transat <- readHelix(transat_file)
> message("RFAM structure in dot bracket format")
> known_file <- system.file("extdata", "vienna.txt", package = "R4RNA")
> known <- readVienna(known_file)
> message("Work with basepairs instead of helices for more flexibility")
> message("Breaks all helices into helices of length 1")
```

```
> transat <- expandHelix(transat)
> known <- expandHelix(known)
```

1.2 Basic Arc Diagram

The standard arc diagram, where the nucleotide sequence is the horizontal line running left to right from 5' to 3' at the bottom of the diagram. Any two bases that base-pair in a secondary structure are connect with an arc.

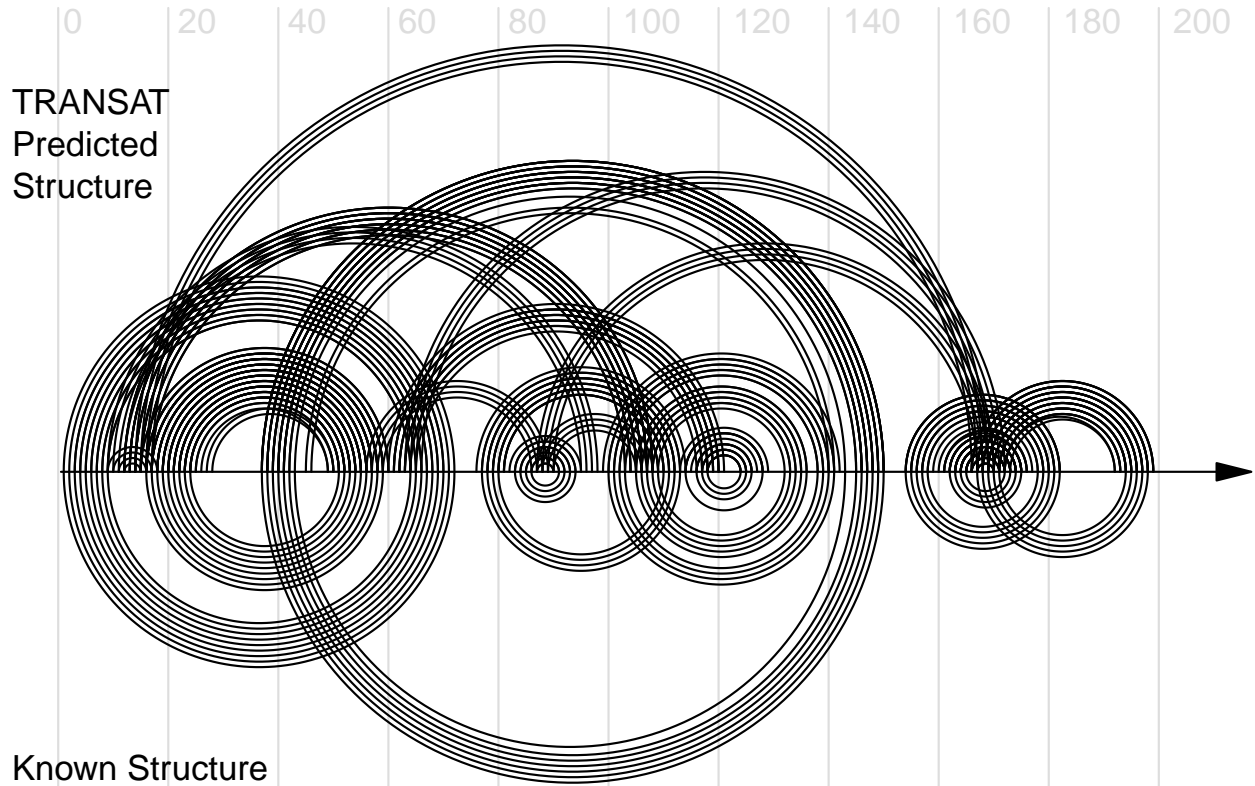
```
> plotHelix(known, line = TRUE, arrow = TRUE)
> mtext("Known Structure", side = 3, line = -2, adj = 0)
```



1.3 Multiple Structures

Two structures for the same sequence can be visualized simultaneously, allowing one to compare and contrast the two structures.

```
> plotDoubleHelix(transat, known, line = TRUE, arrow = TRUE)
> mtext("TRANSAT\nPredicted\nStructure", side = 3, line = -5, adj = 0)
> mtext("Known Structure", side = 1, line = -2, adj = 0)
```



1.4 Filtering Helices

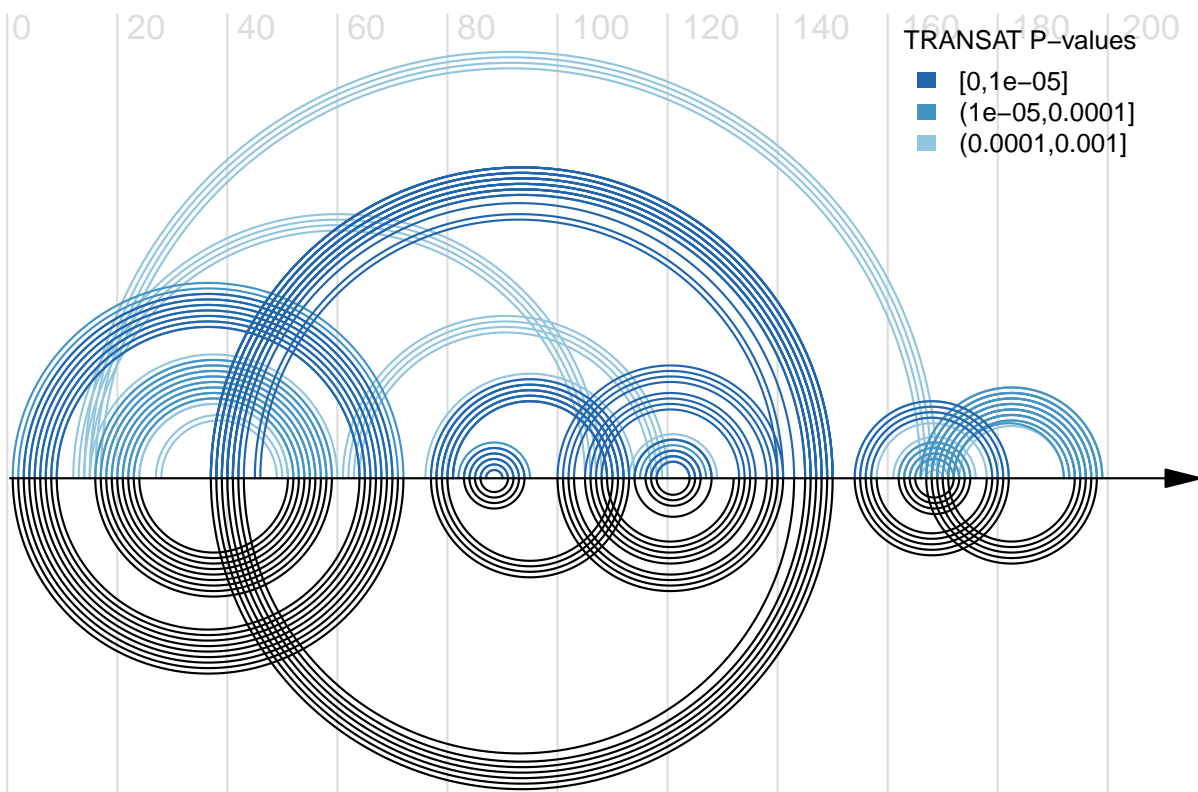
Base-pairs can be associated with a value, such as energy stability or statistical probability, and we can easily filter out basepairs according to such rules.

```
> message("Filter out helices above a certain p-value")
> transat <- transat[which(transat$value <= 1e-3), ]
```

1.5 Colouring Structures

We can also assign colour to the structure according to base-pairs values.

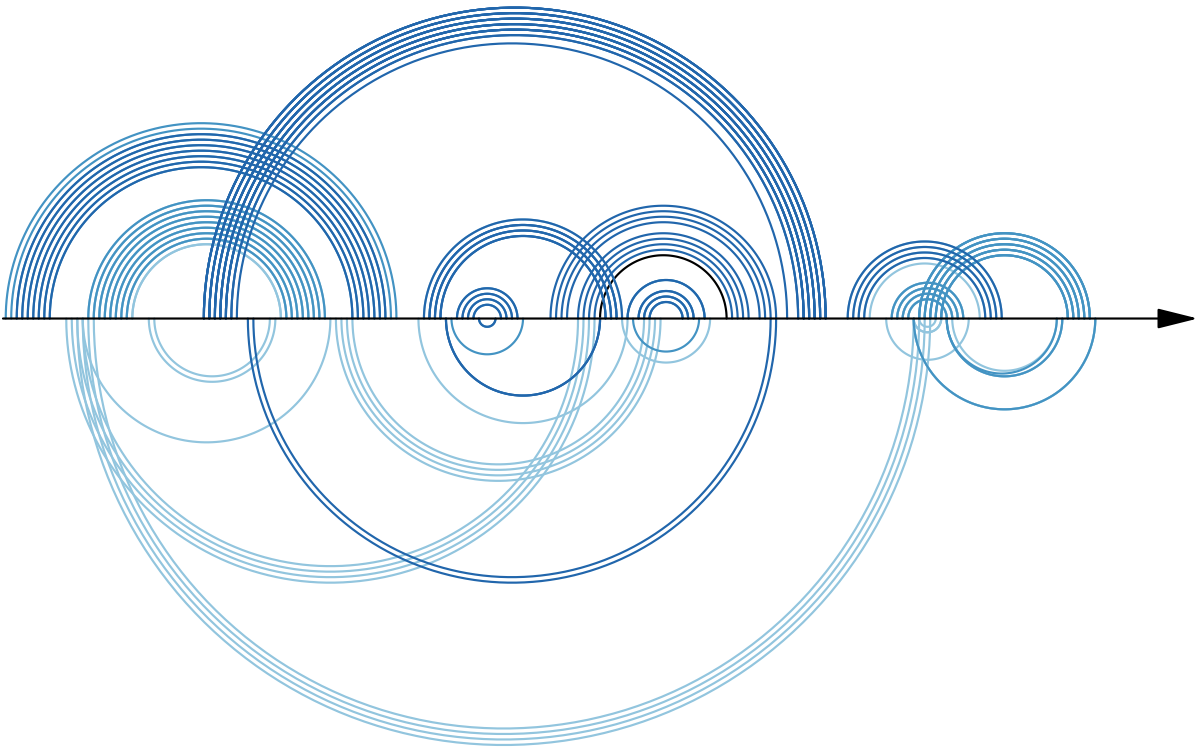
```
> message("Assign colour to basepairs according to p-value")
> transat$col <- col <- colourByValue(transat, log = TRUE)
> message("Coloured encoded in 'col' column of transat structure")
> plotDoubleHelix(transat, known, line = TRUE, arrow = TRUE)
> legend("topright", legend = attr(col, "legend"), fill = attr(col, "fill"),
+       inset = 0.05, bty = "n", border = NA, cex = 0.75, title = "TRANSAT P-values")
```



1.6 Overlapping Multiple Structures

A neat way of visualizing the concordance between two structures is an overlapping structure diagram, which we can use to overlap the predicted TRANSAT structure and the known RFAM structure. Predicted basepairs that exist in the known structure are drawn above the line, and those predicted that are not known to exist are drawn below. Those known but unpredicted are shown in black above the line.

```
> plotOverlapHelix(transat, known, line = TRUE, arrow = TRUE, scale = FALSE)
```

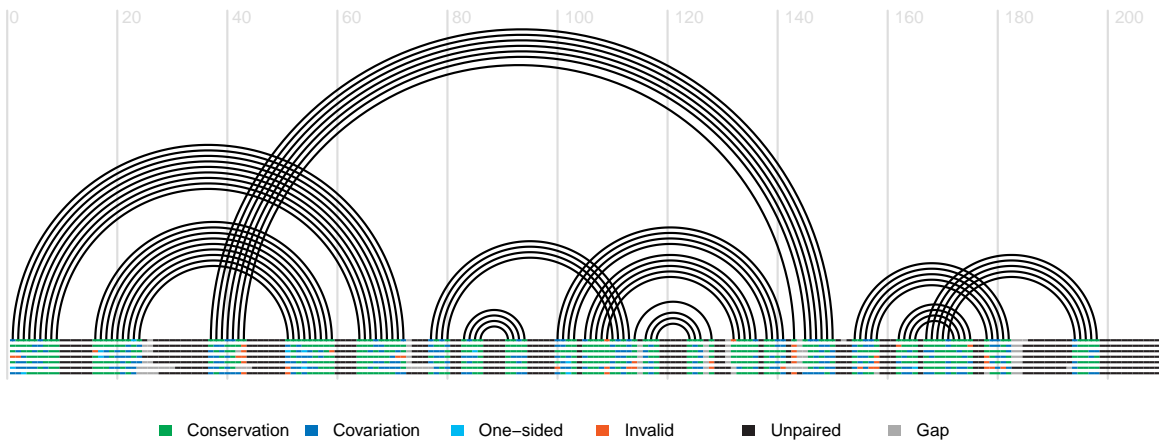


1.7 Visualizing Multiple Sequence Alignments

In addition to visualizing the structure alone, we can also visualize a secondary structure along with aligned nucleotide sequences. In the following, we will read in a multiple sequence alignment obtained from RFAM, and visualize the known structure on top of it.

We can also annotate the alignment colours according to their agreement with the known structure. If a sequence can form as basepair as dictated by the structure, the basepair is coloured green, else red. For green basepairs, if a mutation has occurred, but basepairing potential is retained, it is coloured in blue (dark for mutations in both bases, light for single-sided mutation). Unpaired bases are in black and gaps are in grey.

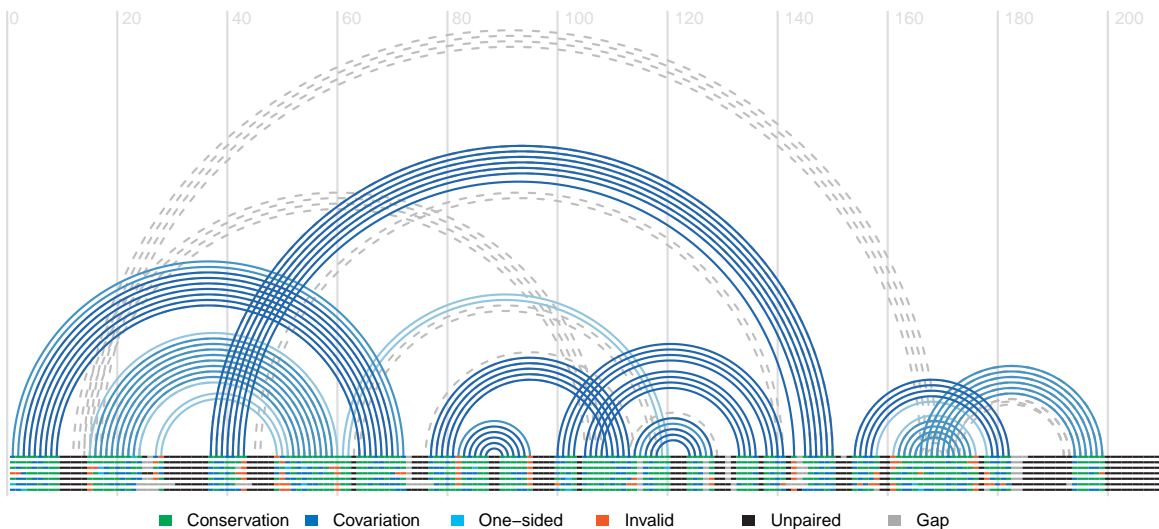
```
> message("Multiple sequence alignment of interest")
> library(Biostrings)
> fasta_file <- system.file("extdata", "fasta.txt", package = "R4RNA")
> fasta <- as.character(readBStringSet(fasta_file))
> message("Plot covariance in alignment")
> plotCovariance(fasta, known, cex = 0.5)
```



1.8 Multiple Sequence Alignments with Annotated Arcs

Arcs can be coloured as usual. It should be noted that structures with conflicting basepairs (arcs sharing a base) cannot be visualized properly on a multiple sequence alignment, and are typically filtered out (*e.g.* drawn in grey here).

```
> plotCovariance(fasta, transat, cex = 0.5, conflict.col = "grey")
```

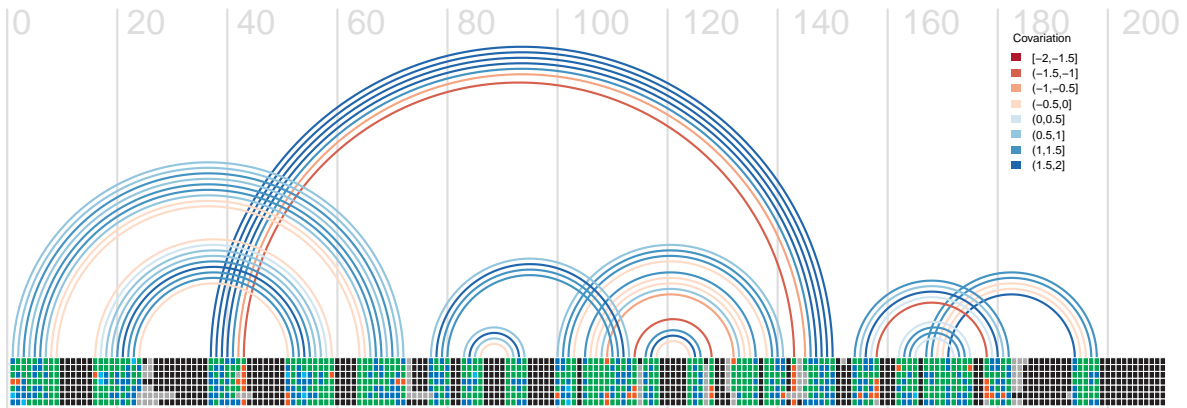


1.9 Additional Colouring Methods

Various other methods of colour arcs exist, along with many options to control appearances:

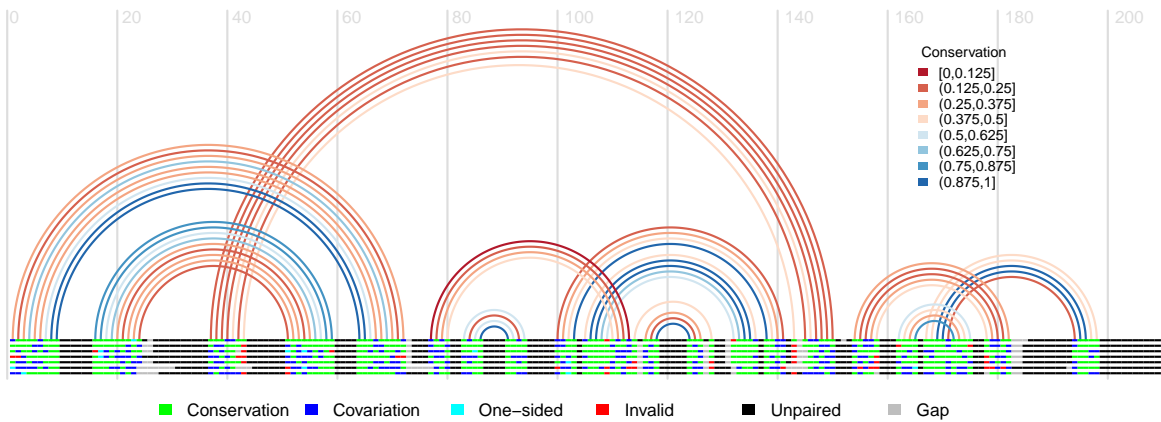
1.9.1 Colour By Covariation (with alignment as blocks)

```
> col <- colourByCovariation(known, fasta, get = TRUE)
> plotCovariance(fasta, col, grid = TRUE, legend = FALSE)
> legend("topright", legend = attr(col, "legend"), fill = attr(col, "fill"),
+       inset = 0.1, bty = "n", border = NA, cex = 0.37, title = "Covariation")
```



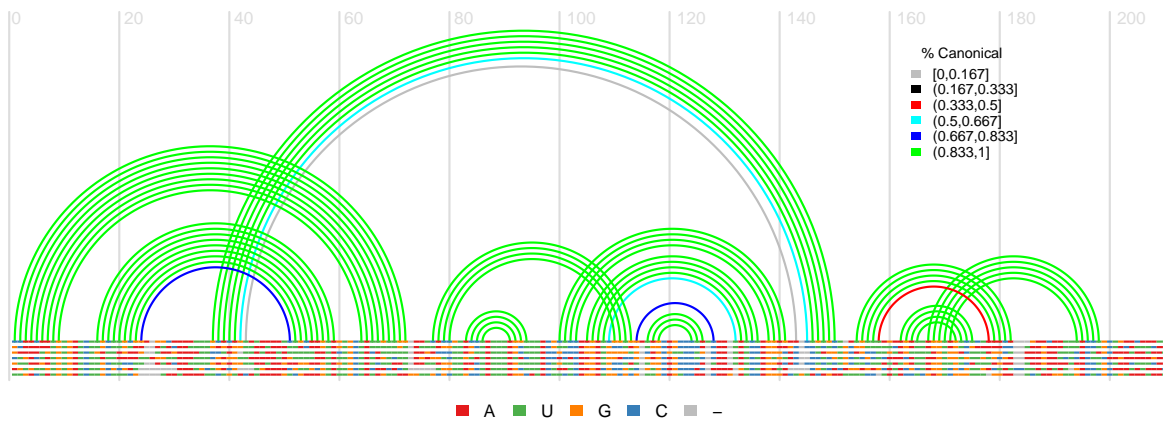
1.9.2 Colour By Conservation (with custom alignment colours)

```
> custom_colours <- c("green", "blue", "cyan", "red", "black", "grey")
> plotCovariance(fasta, col <- colourByConservation(known, fasta, get = TRUE),
+   palette = custom_colours, cex = 0.5)
> legend("topright", legend = attr(col, "legend"), fill = attr(col, "fill"),
+   inset = 0.15, bty = "n", border = NA, cex = 0.75, title = "Conservation")
```



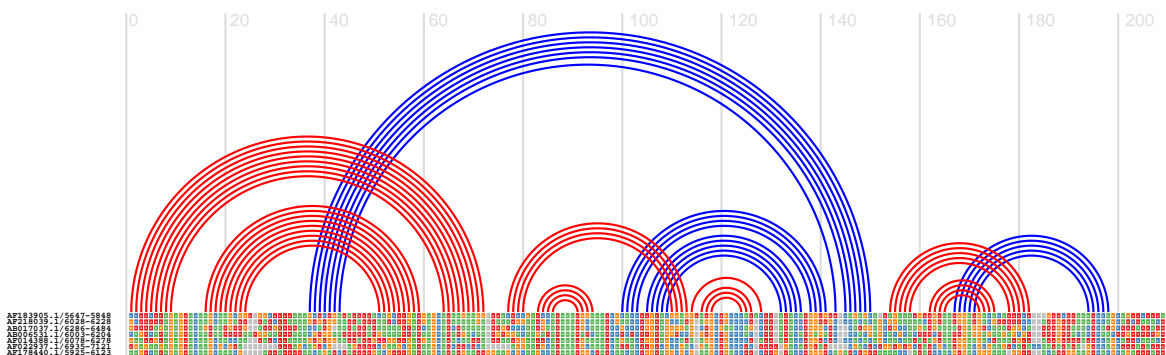
1.9.3 Colour By Percentage Canonical Basepairs (with custom arc colours)

```
> col <- colourByCanonical(known, fasta, custom_colours, get = TRUE)
> plotCovariance(fasta, col, base.colour = TRUE, cex = 0.5)
> legend("topright", legend = attr(col, "legend"), fill = attr(col, "fill"),
+   inset = 0.15, bty = "n", border = NA, cex = 0.75, title = "% Canonical")
```



1.9.4 Colour Pseudoknots (with CLUSTALX-style alignment)

```
> col <- colourByUnknottedGroups(known, c("red", "blue"), get = TRUE)
> plotCovariance(fasta, col, base.colour = TRUE, legend = FALSE, species = 23, grid = TRUE, text = TRUE)
```



2 Session Information

The version number of R and packages loaded for generating the vignette were:

- R version 4.4.1 (2024-06-14), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_GB, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Time zone: America/New_York
- TZcode source: system (glibc)
- Running under: Ubuntu 24.04.1 LTS
- Matrix products: default

- BLAS: `/home/biocbuild/bbs-3.20-bioc/R/lib/libRblas.so`
- LAPACK: `/usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.12.0`
- Base packages: `base`, `datasets`, `grDevices`, `graphics`, `methods`, `stats`, `stats4`, `utils`
- Other packages: `BiocGenerics 0.52.0`, `Biostrings 2.74.0`, `GenomeInfoDb 1.42.0`, `IRanges 2.40.0`, `R4RNA 1.34.0`, `S4Vectors 0.44.0`, `XVector 0.46.0`
- Loaded via a namespace (and not attached): `GenomeInfoDbData 1.2.13`, `R6 2.5.1`, `UCSC.utils 1.2.0`, `compiler 4.4.1`, `crayon 1.5.3`, `httr 1.4.7`, `jsonlite 1.8.9`, `tools 4.4.1`, `zlibbioc 1.52.0`