

AIRR data analysis with immunarch

Bioinformatics software problems and
how to make immune data analysis effortless

AIRR-seq Tool Demonstration | Genoa, Italy | May 2019

Vadim I. Nazarov

National Research University "Higher School of Economics", Russia
ImmunoMind, USA

Table of contents

- 1 Two types of software**
- 2 We aimed for the living**
- 3 What is "painless AIRR analysis"?**
- 4 AIRR data analysis as an one-step process**
- 5 Parsing of AIRR data**
- 6 Visualisation of AIRR data**
- 7 Pipeline example**
- 8 Conclusion & contacts**

Two types of software

Living software

Software for People

- thoroughly developed and tested software that solves a real-world task
- addresses the needs of people beyond the laboratory
- regularly updated and supported
- ideal case does not exist, but we all seek for it

Undead software

Software for Software

- in-house scripts re-branded as a trendy new software package
- yet they are often but a combination of lazy workarounds
- in majority of cases the support and updates cease immediately after publication
- due to these facts the software rarely works as expected if works at all

We aimed for the living

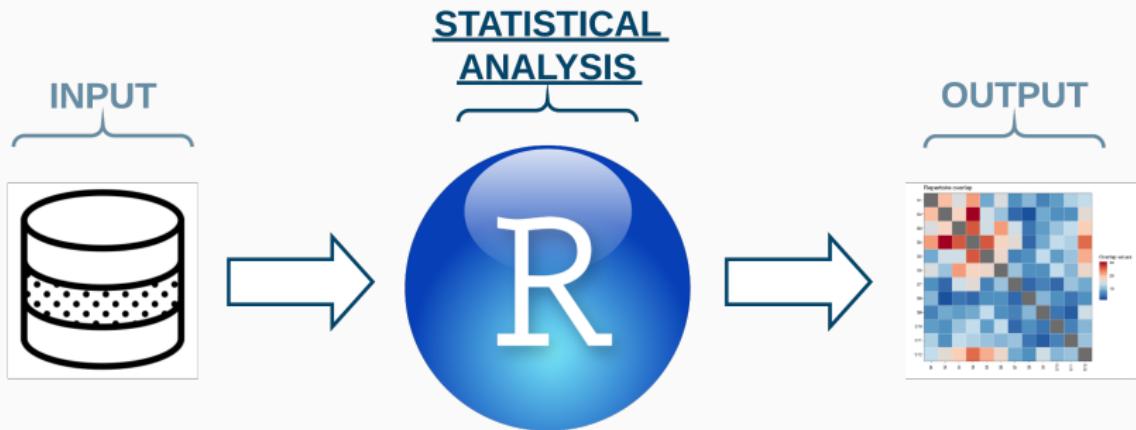
2015 – **tcR**: an R package for AIRR data analysis,
18,000 downloads over 4 years

2017 – we started to build **immunarch** focusing on
painless AIRR data analysis

2018 – start of interviews and surveys with
researchers from different laboratories

2018–2019 – we revealed a couple of our **over-
sights...**

What is "painless AIRR analysis"?



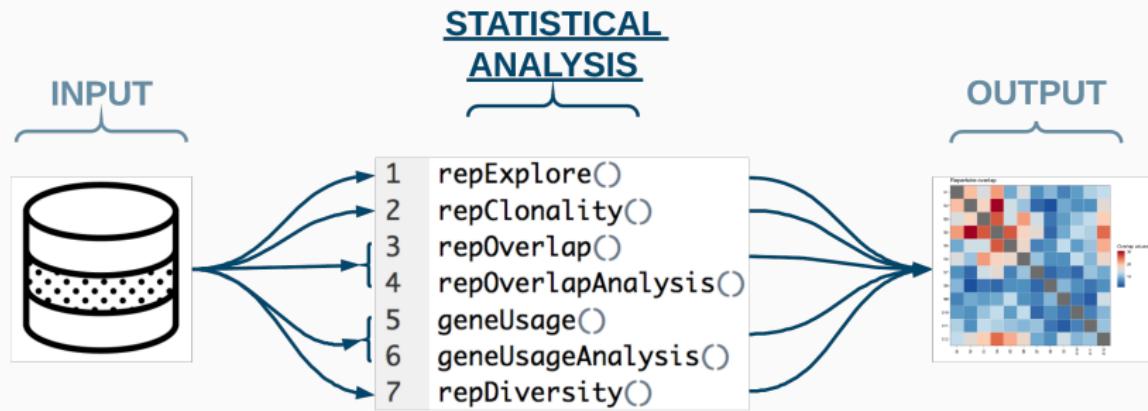
Initial assumption: painless analysis is a fast and effortless computation of common statistics

AIRR data analysis as an one-step process I

We focused on the statistical analysis of AIRR data first and made it as effortless as possible:

```
1 # Overlap of repertoires using Jaccard and Morisita-Horn's indices
2 result = repOverlap(DATA, "jaccard")
3 result = repOverlap(DATA, "morisita")
4
5 # Jensen-Shannon divergence applied to gene usage
6 result = geneUsage(DATA, "HomoSapiens.TRBV")
7 result2 = geneUsageAnalysis(result, "js")
8
9 # Diversity estimation using D50 and Chao1 indices
10 result = repDiversity(DATA, "d50")
11 result = repDiversity(DATA, "chao1")
```

AIRR data analysis as an one-step process II



Why do scientists suffer when they analyse AIRR data?

Problem I: Researchers struggle to **load the multiple data formats** into R

And struggle greatly

AIRR data analysis - parsing II

repLoad() - one function to parse them all. It automatically detects the input format and parses it without any specification from the user

Folder with AIRR data
in different formats

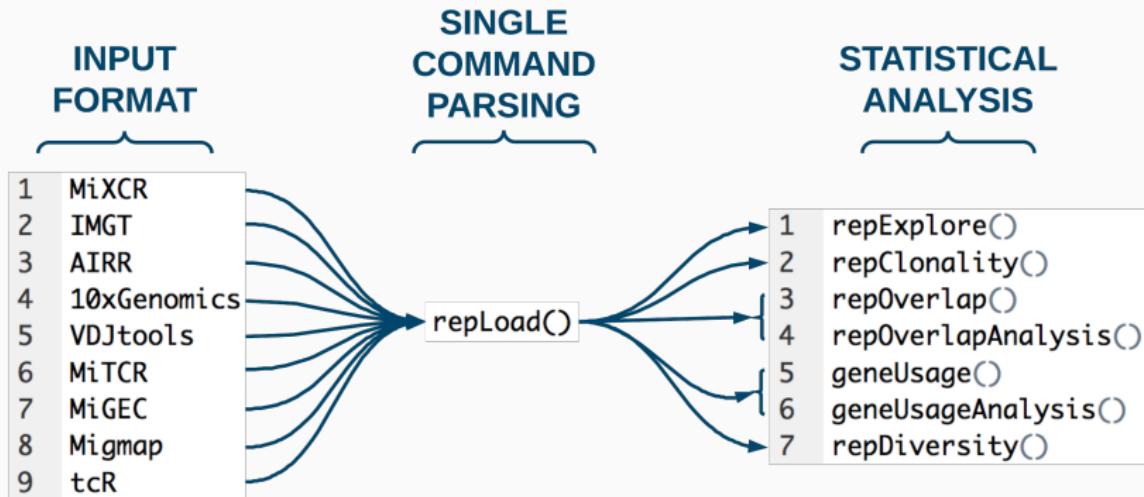
experiment data
imgt.txt
immunoseq_v1.txt
immunoseq_v2.txt
immunoseq_v3.txt
migec.txt
migmap.txt
mitcr.txt
mixcr_old.txt
mixcr_v17.txt
mixcr_v21.txt
tcr.txt
vdjtools_v1.txt
vdjtools_v2.txt



R console: *repLoad()* in action

```
> experiment_data = repLoad("experiment data/")
Parsing experiment data/ ...
Parsing experiment data//imgt.txt -- imgt
Parsing experiment data//immunoseq_v1.txt -- immunoseq
Parsing experiment data//immunoseq_v2.txt -- immunoseq
Parsing experiment data//immunoseq_v3.txt -- immunoseq
Parsing experiment data//migec.txt -- migec
Parsing experiment data//migmap.txt -- migmap
Parsing experiment data//mitcr.txt -- mitcr
Parsing experiment data//mixcr_old.txt -- mixcr
Parsing experiment data//mixcr_v17.txt -- mixcr
Parsing experiment data//mixcr_v21.txt -- mixcr
Parsing experiment data//tcr.txt -- tcr
Parsing experiment data//vdjtools_v1.txt -- vdjtools
Parsing experiment data//vdjtools_v2.txt -- vdjtools
```

AIRR data analysis - the parsing pipeline

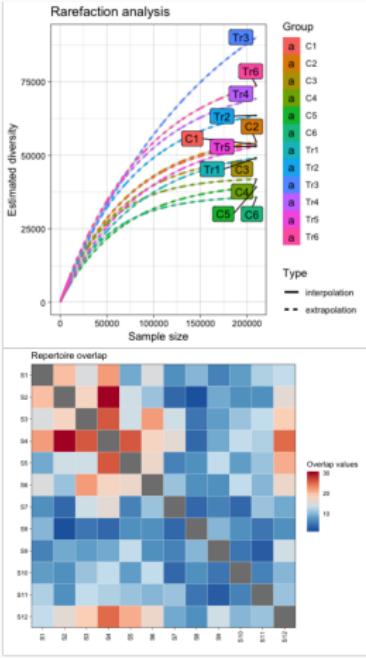


Problem II: Researchers **constantly need to manipulate visualisations** to make them publication-ready, so we designed tools for easy visualisation

Visualisation of AIRR data analysis II

`vis()` - one function to plot them all

```
repDiversity(DATA$data, "raref") %>% vis()
```



Visualisation of AIRR data analysis III

`fixVis()` is a tool to manipulate plots painlessly - change font sizes, titles, angles, etc.

FixVis: make your plots publication-ready already!

Save Plot to R console

General Title & subtitle

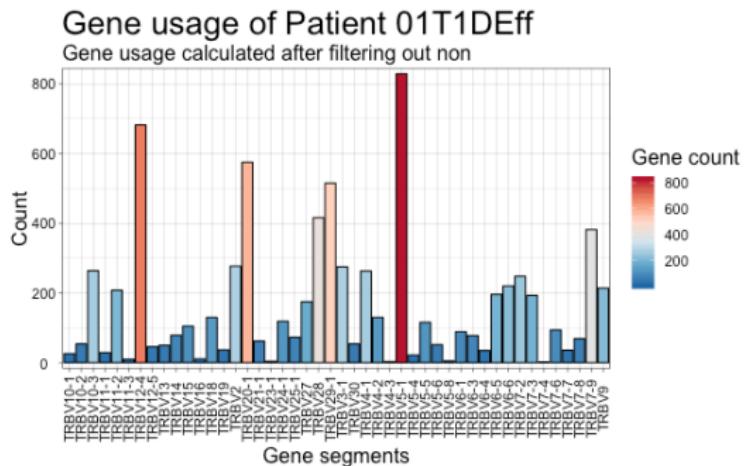
Legends X axis Y axis

Title Subtitle

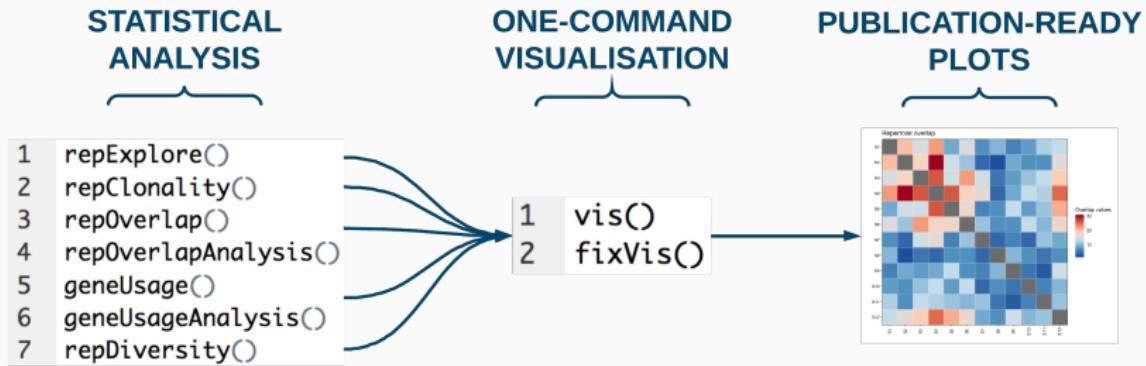
Subtitle text:
Gene usage calculated after filtering out non

Subtitle text size: 16

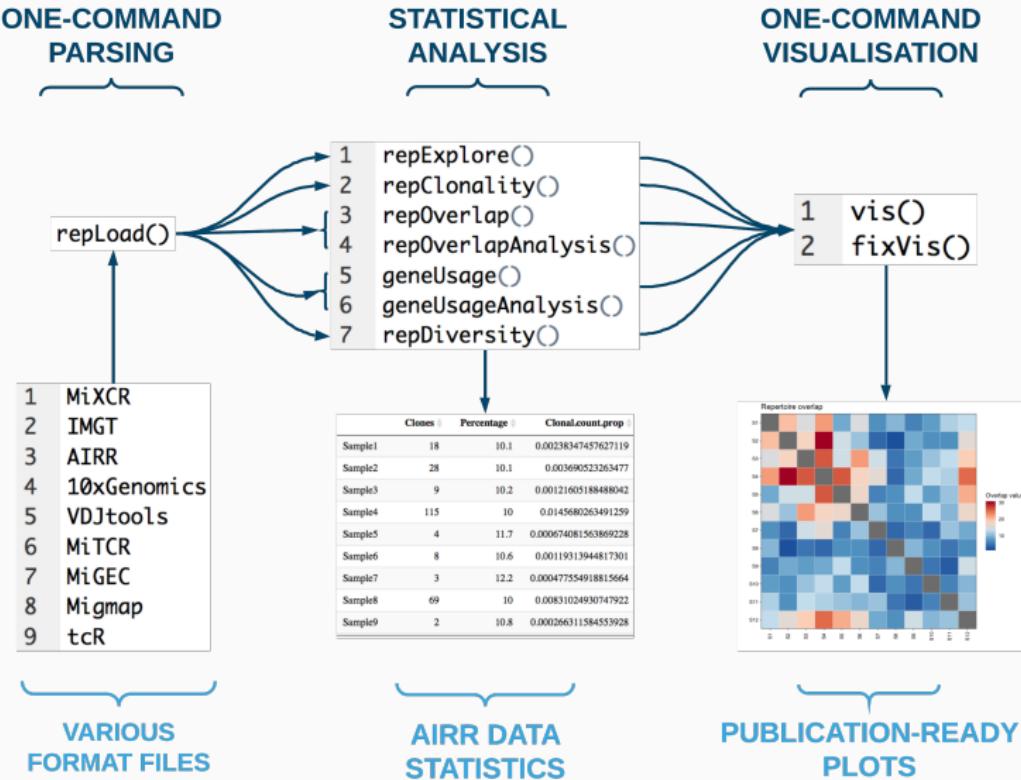
Subtitle text horizontal adjustment:



AIRR data analysis - visualisation pipeline



Overall immunarch pipeline



Pipeline example I

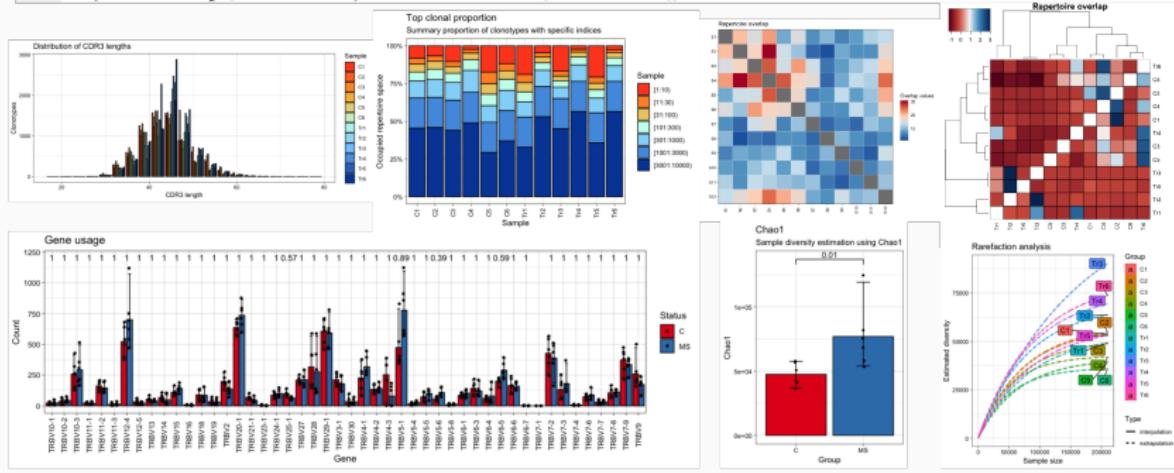
What does it take to perform a typical pipeline in **immunarch**?

1. Load all the data into R
2. Compute and visualise:
 - 2.1 distribution of CDR3 lengths
 - 2.2 relative abundances of the most frequent clonotypes
 - 2.3 the number of public clonotypes between samples
 - 2.4 Morisita-Horn index values between samples
 - 2.5 V gene usage and group the visualisation by treatment status
 - 2.6 Chao1 diversity metric and group the visualisation by treatment status
 - 2.7 rarefaction curves

Pipeline example II

Answer: 8 lines of code and 2 minutes of coding

```
1 DATA = repLoad("path/to/your/folder")
2     repExplore(DATA$data, "len") %>% vis()
3     repClonality(DATA$data, "top") %>% vis()
4         repOverlap(DATA$data) %>% vis()
5     repOverlap(DATA$data, "morisita") %>% vis(.plot="heatmap2")
6         geneUsage(DATA$data) %>% vis(.by="Status", .meta=DATA$meta, .grid=F)
7     repDiversity(DATA$data, "chao1") %>% vis(.by="Status", .meta=DATA$meta)
8     repDiversity(DATA$data, "rarefaction") %>% vis()
```



Thank you!

ImmunoMind

- Evgenii Ofitserov
- Sergey Fedyushchenko
- Vadim I. Nazarov

Collaborators

- Vasily Tsvetkov
- Daniel J. Moore, PhD
- Victor Greiff, PhD

Interviewees and surveyees

- Jean-Philippe Bürkert
- Anne Eugster
- Wyatt J. McDonnell
- Anna Lorenc
- Nils-Petter Rudqvist
- Anna Vyborova
- Alexander G. Watson
- Cindy L. Zuleger
- ... and many more

Conclusion & contacts

Want to focus on research and not coding?

Pre-release version is on
<https://immunarch.com>

Presentation and news:

Twitter @immunomind

Interested in autoimmunity or oncology?

We are looking for partners on biomarker discovery—CDR3 motifs



Vadim I. Nazarov

HSE (RU), ImmunoMind (US)

vadim@immunomind.io

... or in Genoa until
May 19