

TCR Analysis – Principles, Workflows, and Software

Ben K. Margetts

UCL ICH
& Great Ormond Street Hospital

TCR Analysis Handover – 26th March

AutoTCR

AutoTCR is our lab's **automated sample processing and analysis pipeline**

The **data structure** it produces forms **the basis** for all of our **figures and analyses**



To date, our lab has generated **> 600 GB** of compressed TCR sequencing data
That is equal to **~ 600,000 word documents or 15,000,000 text files**

TCR & Data Science/Biology



Without the 'data science' angle of this work, we would have **no results**. Continued development of this is **essential** to being a functional TCR lab.

Lab Biologist – Generate the data
Data Biologist – Make the data meaningful

TCR & Data Science/Biology



Without the 'data science' angle of this work, we would have **no results**. Continued development of this is **essential** to being a functional TCR lab.

Lab Biologist – Generate the data

Lab Biologist

Generate
Data

The way we used to do science

Lab Biologist

Analyse &
Interpret
Data

(The fun part of science!)

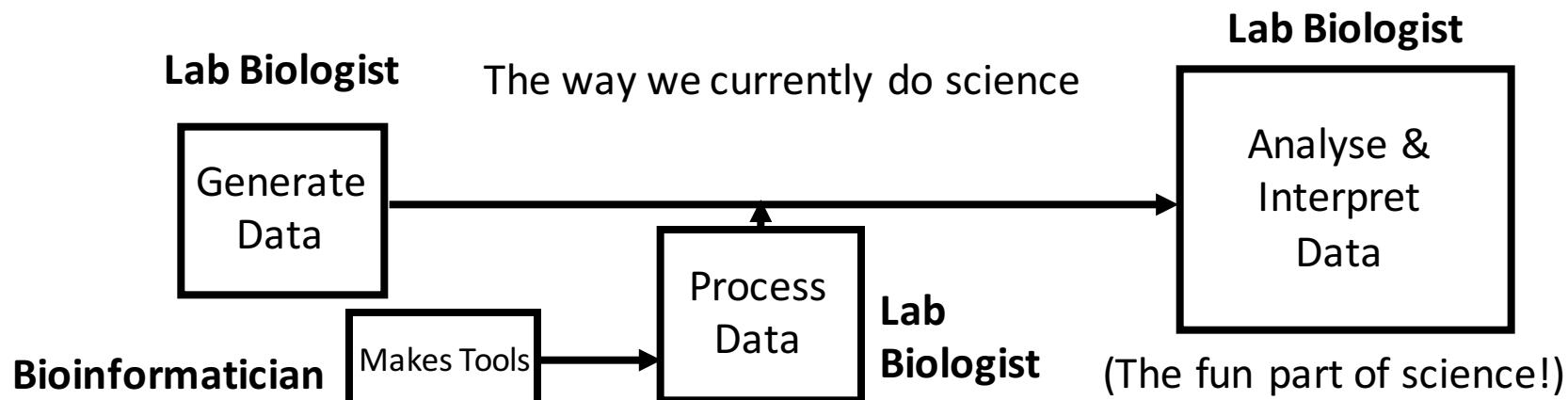
TCR & Data Science/Biology



Without the 'data science' angle of this work, we would have **no results**. Continued development of this is **essential** to being a functional TCR lab.

Lab Biologist – Generate the data

Bioinformatician – Develops methods and tools for Biologists



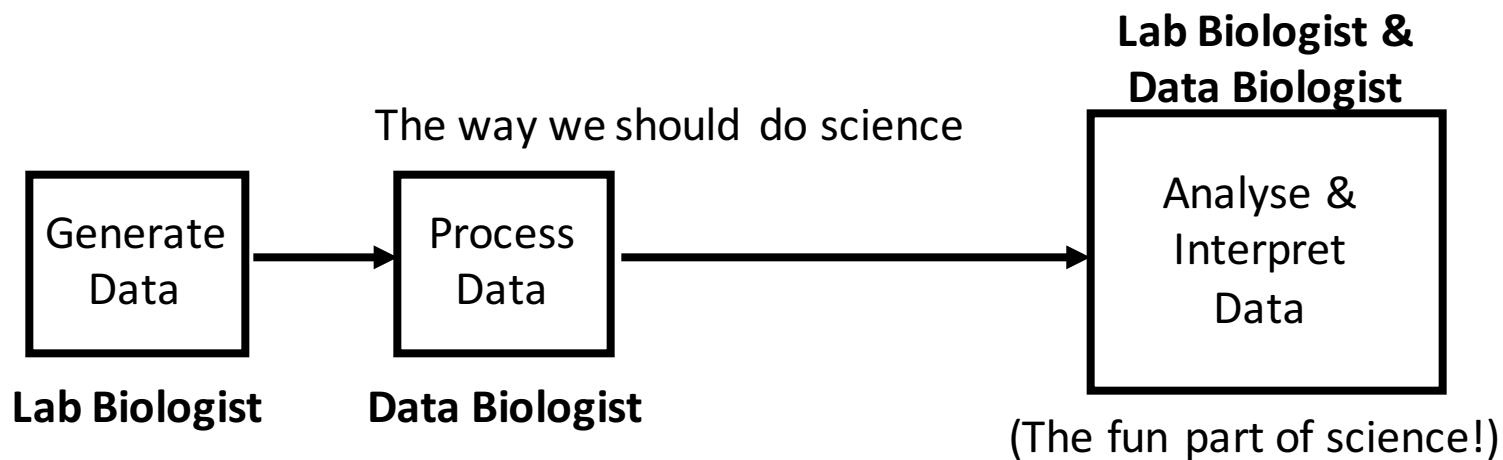
TCR & Data Science/Biology



Without the 'data science' angle of this work, we would have **no results**. Continued development of this is **essential** to being a functional TCR lab.

Lab Biologist – Generate the data

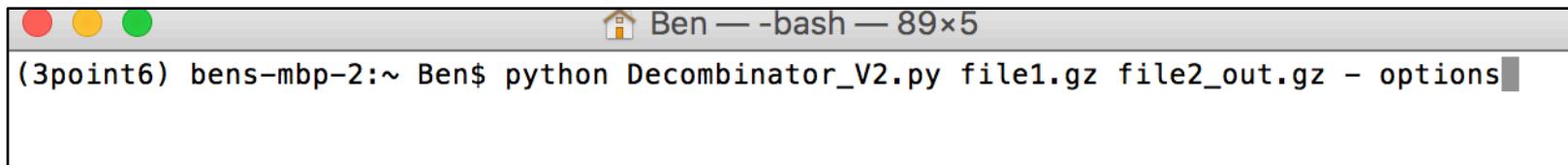
Data Biologist – Make the data meaningful



Why Make The Script?

Before (est. time = days):

- Manually run each file through each of the 5 Decombinator Scripts (12 files per run)



A screenshot of a macOS terminal window. The title bar says "Ben — -bash — 89x5". The command line shows: "(3point6) bens-mbp-2:~ Ben\$ python Decombinator_V2.py file1.gz file2_out.gz - options".

- Validate that no mistakes have been made
- Read in data to analysis platform (Excel/R/Python/SAS/STATA)
- Assign metadata to each file (time, patient, disease, etc)
- Structure data into a dataframe
- Produce data visualisations

Now (est. time = 3 ½ hours):

- Provide correct name in index file.
- Leave AutoTCR running for ~ 3 hours.
- Done.

Unfortunately,

Hadley Wickham – *'You can't do data science in a GUI'*

	Python	R	R Studio & Tidyverse	Bash
1. Installation & Requirements				
2. Automation and File Structure				
3. Reproducibility				
4. Decombinator				
a) Demultiplexor				
b) Decombinator				
c) CDR3Translator				
d) GNTranslator				
5. Building a Dataframe				
6. Automated Analysis				
7. Debugging/Common Issues				
8. Subsampling				

Installation & Requirements

- Python Version 2.7

(Anaconda distribution <https://www.anaconda.com/download/#macos>)

- Python packages (using Conda package manager)

(see <https://github.com/innate2adaptive/Decombinator#installation>)

- R

(see <https://cloud.r-project.org>)

- R Studio

(see <https://www.rstudio.com/products/rstudio/download/>)

- R packages (installed by AutoTCR script)

- Mac OSX only - command line tools

(type *xcode-select –install* in Terminal)

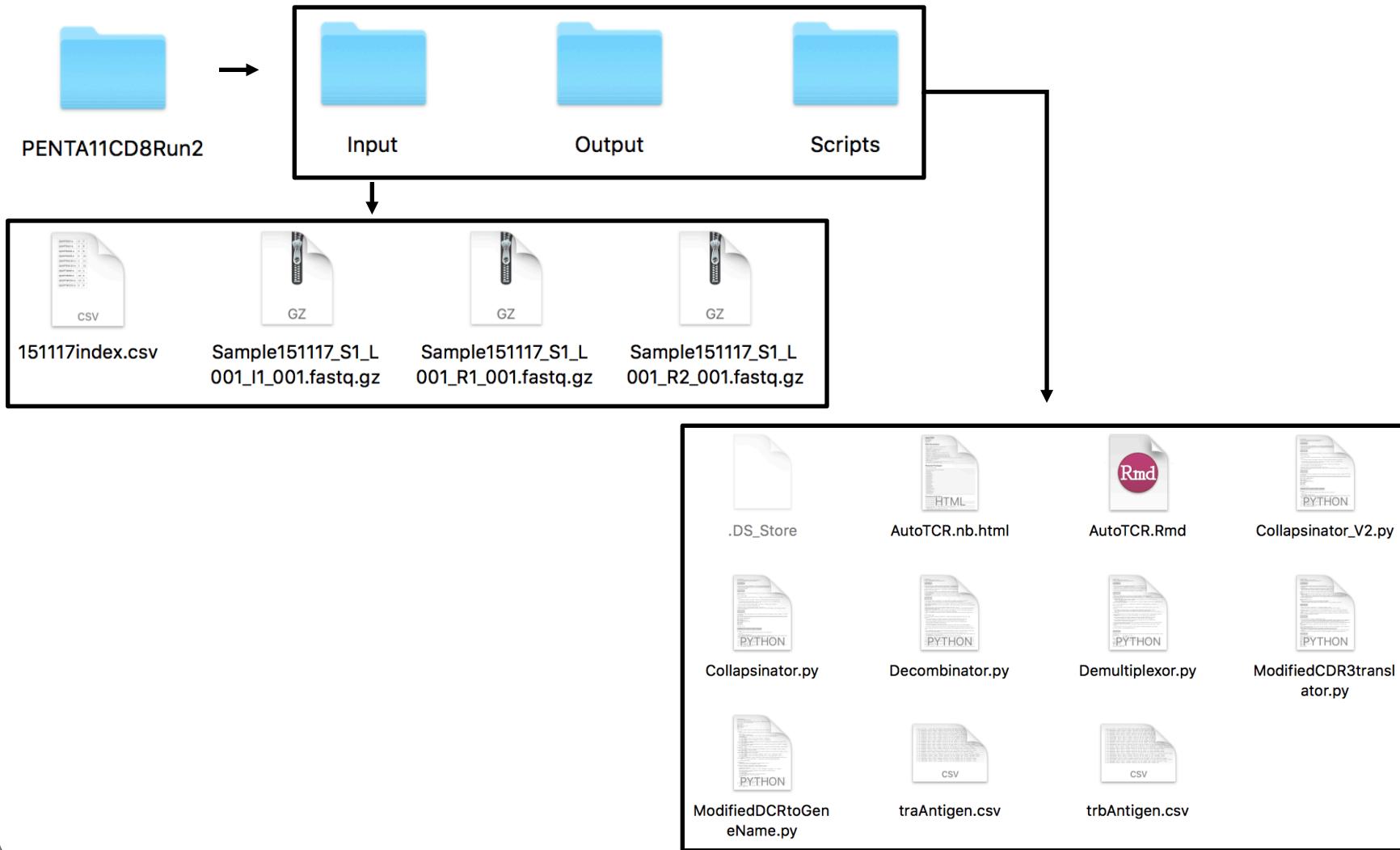
- If in doubt, please ask.

Automation & File Structure

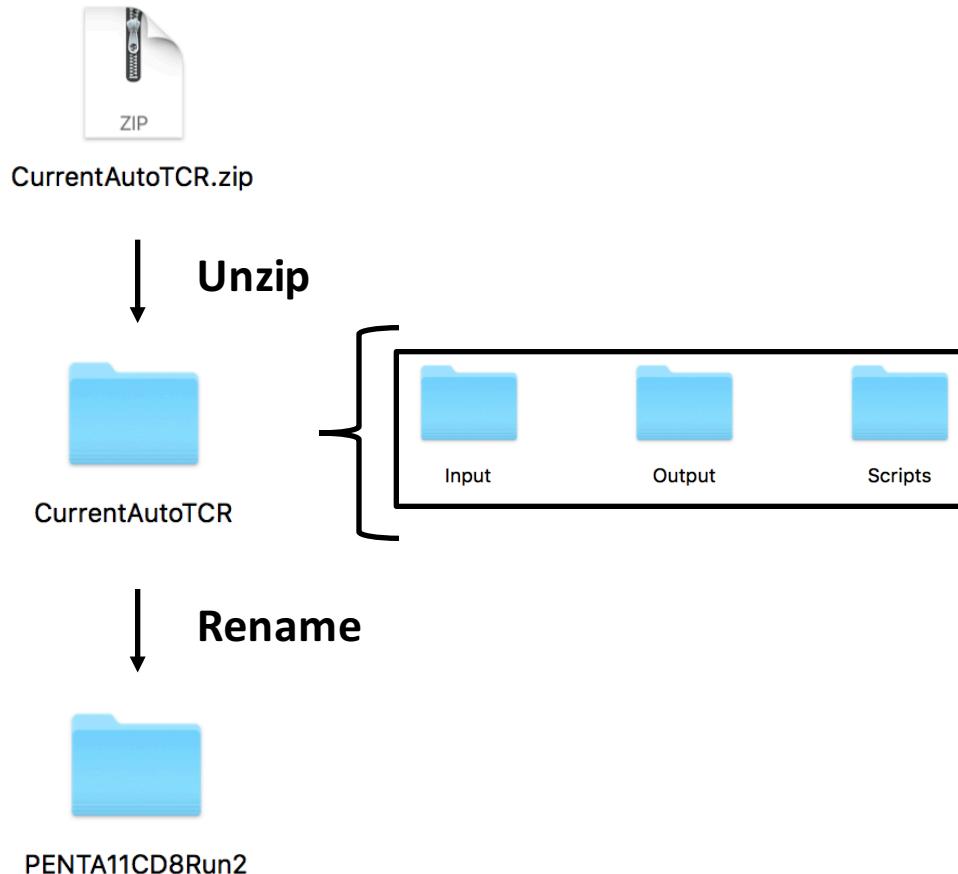
What do you need to do to run AutoTCR?

1. Create your file structure from template (Input, Output, Scripts)
2. Put your fastq files in ‘Input’
3. Create your index file (and put it in ‘Input’)
4. Open ‘AutoTCR.Rmd’ and edit the first section

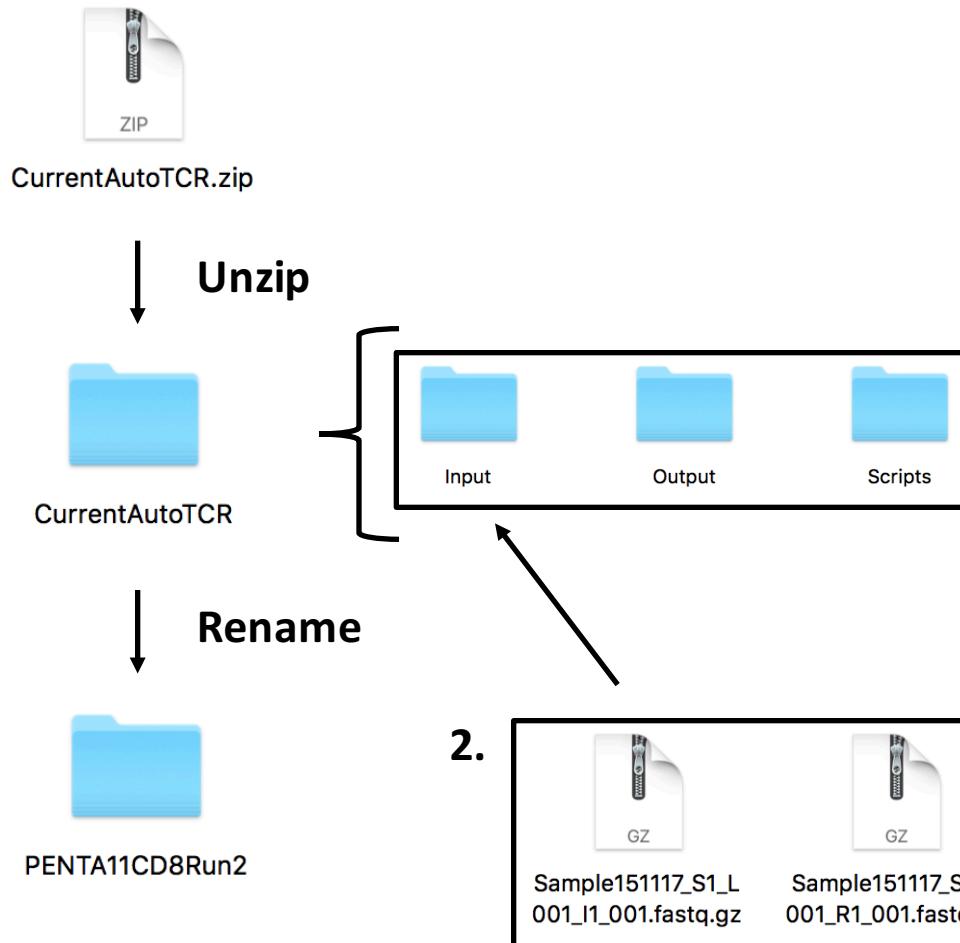
Automation & File Structure



1. Create File Structure From Template



1. Create File Structure From Template



3. The Index File

Sample Name	Index	
QUFCTW0-a	11	1
QUFCTW0-b	12	2
QUDCTW150-a	13	3
QUDCTW150-b	5	4
QUTPTIW48-a	4	5
QUTPTIW48-b	2	6
QUQPTIW48-a	3	7
QUQCTW150-a	2	11
QUQCTW150-b	3	12

→ ‘081217index.csv’

- ‘.csv’ – comma separated value

- Create with a text editor (i.e. bbedit)

1	QUFCTW0-a,11,1
2	QUFCTW0-b,12,2
3	QUDCTW150-a,13,3
4	QUDCTW150-b,5,4
5	QUTPTIW48-a,4,5
6	QUTPTIW48-b,2,6
7	QUQPTIW48-a,3,7
8	QUQCTW150-a,2,11
9	QUQCTW150-b,3,12

QUFCTW0-a

3. The Index File

Sample Name	Index	
QUFCTW0-a	11	1
QUFCTW0-b	12	2
QUDCTW150-a	13	3
QUDCTW150-b	5	4
QUTPTIW48-a	4	5
QUTPTIW48-b	2	6
QUQPTIW48-a	3	7
QUQCTW150-a	2	11
QUQCTW150-b	3	12

→ ‘081217index.csv’

- ‘.csv’ – comma separated value

- Create with a text editor (i.e. bbedit)

1	QUFCTW0-a,11,1
2	QUFCTW0-b,12,2
3	QUDCTW150-a,13,3
4	QUDCTW150-b,5,4
5	QUTPTIW48-a,4,5
6	QUTPTIW48-b,2,6
7	QUQPTIW48-a,3,7
8	QUQCTW150-a,2,11
9	QUQCTW150-b,3,12

QUFCTW0-a

QUFCTW0-a - Patient name

3. The Index File

Sample Name	Index	
QUFCTW0-a	11	1
QUFCTW0-b	12	2
QUDCTW150-a	13	3
QUDCTW150-b	5	4
QUTPTIW48-a	4	5
QUTPTIW48-b	2	6
QUQPTIW48-a	3	7
QUQCTW150-a	2	11
QUQCTW150-b	3	12

→ ‘081217index.csv’

- ‘.csv’ – comma separated value

- Create with a text editor (i.e. bbedit)

1	QUFCTW0-a,11,1
2	QUFCTW0-b,12,2
3	QUDCTW150-a,13,3
4	QUDCTW150-b,5,4
5	QUTPTIW48-a,4,5
6	QUTPTIW48-b,2,6
7	QUQPTIW48-a,3,7
8	QUQCTW150-a,2,11
9	QUQCTW150-b,3,12

QUFCTW0-a

QUFCTW0-a - Patient name

QUFCTW0-a - Time value (unit is irrelevant)

3. The Index File

Sample Name	Index	
QUFCTW0-a	11	1
QUFCTW0-b	12	2
QUDCTW150-a	13	3
QUDCTW150-b	5	4
QUTPTIW48-a	4	5
QUTPTIW48-b	2	6
QUQPTIW48-a	3	7
QUQCTW150-a	2	11
QUQCTW150-b	3	12

→ ‘081217index.csv’

- ‘.csv’ – comma separated value

- Create with a text editor (i.e. bbedit)

1	QUFCTW0-a,11,1
2	QUFCTW0-b,12,2
3	QUDCTW150-a,13,3
4	QUDCTW150-b,5,4
5	QUTPTIW48-a,4,5
6	QUTPTIW48-b,2,6
7	QUQPTIW48-a,3,7
8	QUQCTW150-a,2,11
9	QUQCTW150-b,3,12

QUFCTW0-a

QUFCTW0-a - Patient name

QUFCTW0-a - Time value (unit is irrelevant)

QUFCTW0-a - Chain (a = alpha, b = beta)

4. Edit AutoTCR

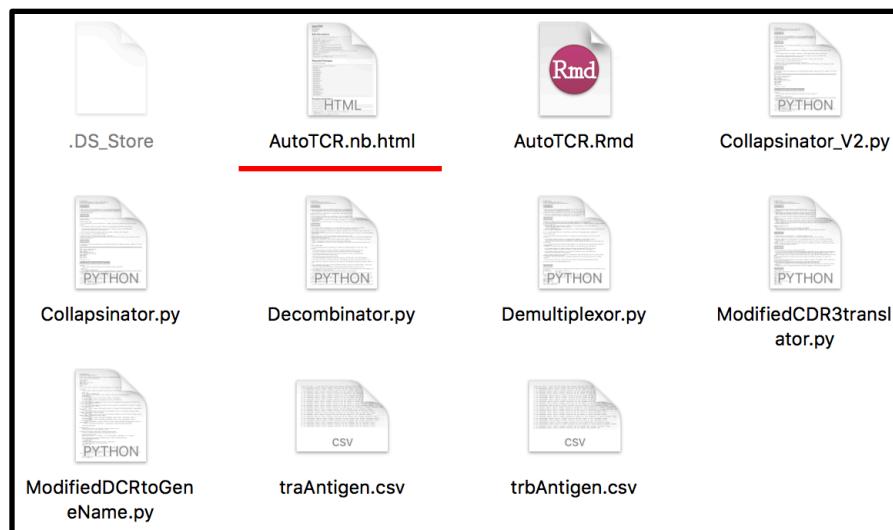
```
1  title: "AutoTCR"
2  author: "Ben Margetts"
3  date: "29/09/2017"
4  output:
5    pdf_document: default
6    html_notebook: default
7    html_document: default
8  ---
9
10 #Edit this section:
11
12 ````{r}
13 #Run ID - what would you like to call this run?
14 runID <- 'pentarun5repeat'
15
16 #Email address to send the run report to
17 emailAddress <- 'ucbpmar@ucl.ac.uk'
18 sendEmail <- T #TRUE or FALSE
19
20 #Contains all data files
21 inputPath <- '~/Desktop/PENTA11Run5Repeat/Input'
22
23 #Path to write output to
24 outputPath <- '~/Desktop/PENTA11Run5Repeat/Output'
25
26 #Contains all decombegrator scripts and this .Rmd script
27 scriptPath <- '~/Desktop/PENTA11Run5Repeat/Scripts'
28
29 #Contains a python installation with all prerequisite decombegrator packages installed
30 python <- '~/anaconda2/bin/python'
31
32 #alpha and beta, or just 1?
33 numberChains <- 2
34
35 #Set to TRUE to run collapsinator with -N, otherwise, set to false.
36 collapsinatorN <- T #TRUE or FALSE
37
38 ````
```

Cmd + S (Save file), Cmd + Alt + R (Run all sections)

Reproducibility

Cmd + Alt + R (Run all sections)

- Produces a .html (web) file documenting the entire run and all output.
- All programmatic analysis is inherently reproducible. If the code is correct, it won't make a mistake



Reproducibility

A screenshot of a Mac OS X desktop environment showing a web browser window. The window title is 'file:///Volumes/Ben-TCR-HDD/Desktop%20Backup%202019%20April%20'. The browser has several tabs open: 'The Health Foundation', 'www.health.org.uk/sites/health/fil...', 'www.health.org.uk/sites/health/f...', 'UCL – Events Calendar – Integr...', and 'AutoTCR'. The main content area displays an R script for 'AutoTCR'.

Code ▾

AutoTCR

Ben Margetts
03/08/2017

Edit this section:

Hide

```
#Run ID - what would you like to call this run?  
runID <- 'test1'  
#Email address to send the run report to  
emailAddress <- 'ucbpmar@ucl.ac.uk'  
sendEmail <- TRUE #TRUE or FALSE  
#Contains all data files  
inputPath <- '~/Google/TCR/AutoTCR/AutoTCRTestV3/Input'  
#Path to write output to  
outputPath <- '~/Google/TCR/AutoTCR/AutoTCRTestV3/Output'  
#Contains all decombinator scripts and this .Rmd script  
scriptPath <- '~/Google/TCR/AutoTCR/AutoTCRTestV3/Scripts'  
#Contains a python installation with all prerequisite decombinator packages installed  
python <- '~/anaconda/bin/python'  
#alpha and beta, or just 1?  
numberChains <- 2  
#Set to TRUE to run collapsinator with -N, otherwise, set to false.  
collapsinatorN <- FALSE #TRUE or FALSE
```

Required Packages

Load in the following packages:

Hide

```
#devtools::install_github("slowkow/ggrepel")  
library(data.table)  
library(zoo)  
library(plyr)  
library(scales)  
library(ggplot2)
```

Decombinator Pipeline

Decombinator: a tool for fast, efficient gene assignment in T-cell receptor sequences using a finite state machine

Niclas Thomas, James Heather, Wilfred Ndifon, John Shawe-Taylor, Benjamin Chain 

Bioinformatics, Volume 29, Issue 5, 1 March 2013, Pages 542–550, <https://doi.org/10.1093/bioinformatics/btt004>

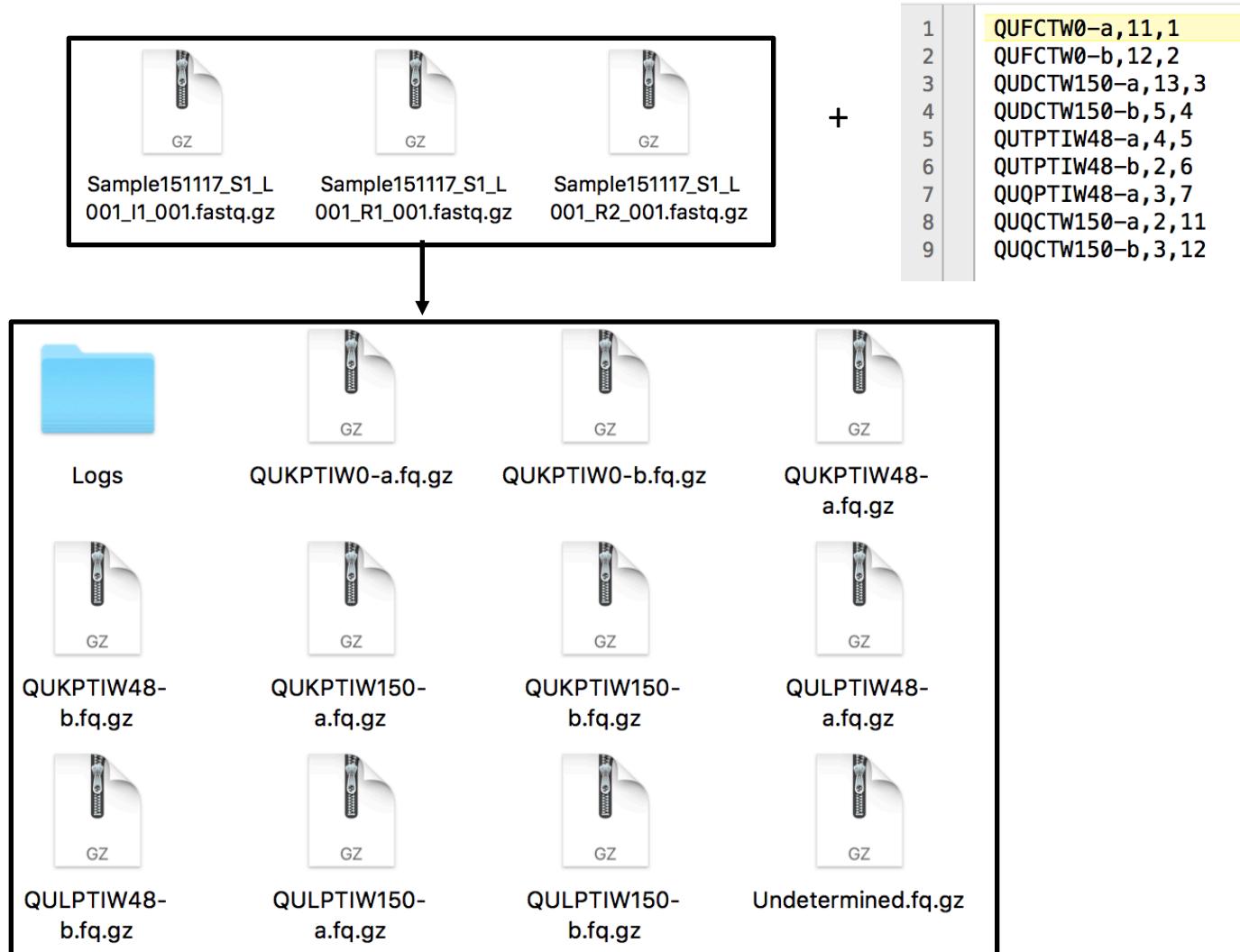
Published: 09 January 2013 Article history ▾

Abstract

Summary: High-throughput sequencing provides an opportunity to analyse the repertoire of antigen-specific receptors with an unprecedented breadth and depth. However, the quantity of raw data produced by this technology requires efficient ways to categorize and store the output for subsequent analysis. To this end, we have defined a simple five-item identifier that uniquely and unambiguously defines each TcR sequence. We then describe a novel application of finite-state automaton to map Illumina short-read sequence data for individual TcRs to their respective identifier. An extension of the standard algorithm is also described, which allows for the presence of single-base pair mismatches arising from sequencing error. The software package, named Decombinator, is tested first on a set of artificial *in silico* sequences and then on a set of published human TcR- β sequences. Decombinator assigned sequences at a rate more than two orders of magnitude faster than that achieved by classical pairwise alignment algorithms, and with a high degree of accuracy (>88%), even after introducing up to 1% error rates in the *in silico* sequences. Analysis of the published sequence dataset highlighted the strong V and J usage bias observed in the human peripheral blood repertoire, which seems to be

1. Demultiplexor
2. Decombinator
3. Collapsinator
4. CDR3 Translator
5. GN Translator

Demultiplexor



Decombinator



Decombinator

Looks for rearranged TCRs

43, 11, 0, 2,
ATCGGGGGAGGGGG,
M01520:11:000000000-
BHDMP:1:1101:13850:1643,

V index, J index, # V deletions, # J
deletions, insert, ID



dcr_beta_QUKPTIW
0-b.n12.gz

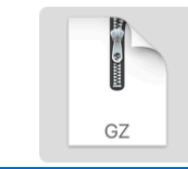
Collapsinator

Productive sequences?

Not PCR amplified?

43, 11, 0, 2,
ATCGGGGGAGGGGG,
M01520:11:000000000-
BHDMP:1:1101:13850:1643,

* 3



dcr_beta_QUKPTIW
0-b.freq.gz

CDR3 Translator & GN Translator

43, 11, 0, 2,
ATCGGGGGAGGGGG,
M01520:11:000000000-
BHDMP:1:1101:13850:1643,

* 3

43, 11,
tgcgcgagcagcttggcgtg
aacagcgattataacctt

* 3

43, 11, CASSFGVNSDYTF

* 3

TRBV13-1, TRBJ-1, CASSFGVNSDYTF

* 3



dcr_beta_QUKPTI
W0-b.cdr3.gz

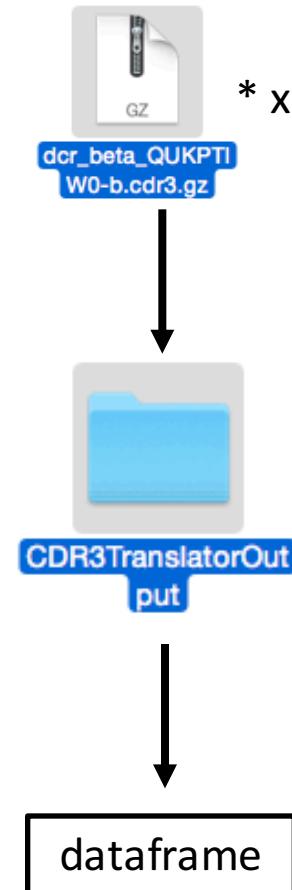


dcr_beta_QUKPTI
W0-b.dcrcdr3.gz

Building a Dataframe

Key elements to capture:

- CDR3 Sequence
- Chain (alpha/beta)
- ID
- Frequency
- Additional computations (for each sequence)



Building a Dataframe



CAAIRGAGSYQLTF	1
CAAIRQAGTALIF	1
CAAIVAGTALIF	1
CAAIVIPGSQGNLIF	1
CAALWGYSASKIIIF	1
CAAMDSNMDNSNYQLIW	1
CAAPHRCGNEKLTF	11
CAAPHRGNEKLTF	209

QUKPTIW0-b = Patient name

+ QUKPTIW0-b = Time value

QUKPTIW0-b = Chain



CAAPTPLYNTDKLF	2
CAARETGASKLTF	1
CAARETSGSRLTF	1
CAARKAETGRGSRLTF	1
CAARKAETSGSRLTF	103
CAARKTSGSRLTF	1
CAARRAETSGSRLTF	1
CAARVRDMRFF	12

QUKPTIW0-a = Patient name

+ QUKPTIW0-a = Time value

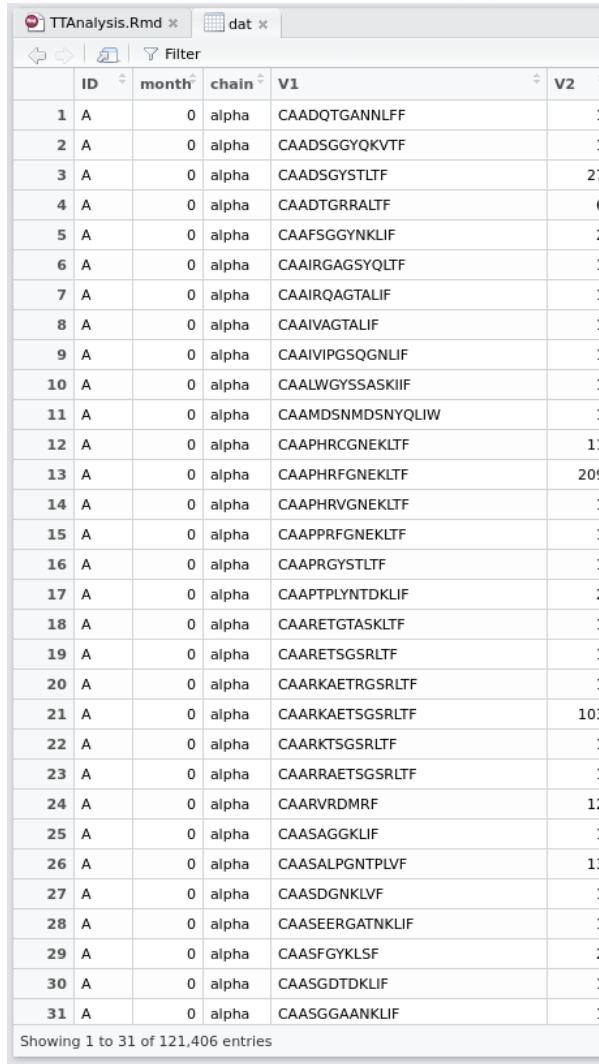
QUKPTIW0-a = Chain

||

dataframe

Building a Dataframe

dataframe



	ID	month	chain	V1	V2
1	A	0	alpha	CAADQTGANNLFF	1
2	A	0	alpha	CAADSGGYQKVTF	1
3	A	0	alpha	CAADSGYSTLTF	27
4	A	0	alpha	CAADTGRRAALTF	6
5	A	0	alpha	CAAFSGGYNKLIF	2
6	A	0	alpha	CAAIRGAGSYQLTF	1
7	A	0	alpha	CAAIRQAGTALIF	1
8	A	0	alpha	CAAIVAGTALIF	1
9	A	0	alpha	CAAIVIPGSQGNLIF	1
10	A	0	alpha	CAALWGYSSASKIIF	1
11	A	0	alpha	CAAMDSNMDNSYQLIW	1
12	A	0	alpha	CAAPHRCGNEKLT	11
13	A	0	alpha	CAAPHRFGNENKLT	209
14	A	0	alpha	CAAPHRVGNEKLT	1
15	A	0	alpha	CAAPPRFGNEKLT	1
16	A	0	alpha	CAAPRGYSTLTF	1
17	A	0	alpha	CAAPTPLYNTDKLIF	2
18	A	0	alpha	CAARETGATSKLT	1
19	A	0	alpha	CAARETSGSRLT	1
20	A	0	alpha	CAARKAETRGSRLT	1
21	A	0	alpha	CAARKAETSGSRLT	103
22	A	0	alpha	CAARKTSGSLT	1
23	A	0	alpha	CAARRAETSGSRLT	1
24	A	0	alpha	CAARVRDMRF	12
25	A	0	alpha	CAASAGGKLIF	1
26	A	0	alpha	CAASALPGNTPLVF	13
27	A	0	alpha	CAASDGNKLVF	1
28	A	0	alpha	CAASEERGATNKLF	1
29	A	0	alpha	CAASFYKLSF	2
30	A	0	alpha	CAASGDTDKLIF	1
31	A	0	alpha	CAASGGAANKLIF	1

Showing 1 to 31 of 121,406 entries

Building a Dataframe dat\$CDR3Len <- length(strsplit(dat\$V1, by = ""))



Building a Dataframe

dataframe



R TTAnalysis.Rmd x dat x

Filter

ID	month	chain	V1	V2	CDR3Len
1	A	0	alpha CAADQTGANNLFF	1	13
2	A	0	alpha CAADSGGYQKVTF	1	13
3	A	0	alpha CAADSGYSTLTF	27	12
4	A	0	alpha CAADTGRRALTF	6	12
5	A	0	alpha CAAFSGGYNKLIF	2	13
6	A	0	alpha CAAIRGAGSYQLTF	1	14
7	A	0	alpha CAAIRQAGTALIF	1	13
8	A	0	alpha CAAIVAGTALIF	1	12
9	A	0	alpha CAAIVIPGSQGNLIF	1	15
10	A	0	alpha CAALWGYSSASKIIIF	1	15
11	A	0	alpha CAAMDSNMDNSNYQLIW	1	16
12	A	0	alpha CAAPHRCGNEKLT	11	14
13	A	0	alpha CAAPHRFGNENKLT	209	14
14	A	0	alpha CAAPHRVGNEKLT	1	14
15	A	0	alpha CAAPPRFGNEKLT	1	14
16	A	0	alpha CAAPRGYSTLTF	1	12
17	A	0	alpha CAAPTPLYNTDKLIF	2	15
18	A	0	alpha CAARETGATSKLT	1	14
19	A	0	alpha CAARETSGSRLT	1	13
20	A	0	alpha CAARKAETRGSRLT	1	15
21	A	0	alpha CAARKAETSGSRLT	103	15
22	A	0	alpha CAARKTGSRLT	1	13
23	A	0	alpha CAARRAETSGSRLT	1	15
24	A	0	alpha CAARVRDMRF	12	10
25	A	0	alpha CAASAGGKLIF	1	11
26	A	0	alpha CAASALPGNTPLVF	13	14
27	A	0	alpha CAASDGNKLVF	1	11
28	A	0	alpha CAASEERGATNKLIF	1	15
29	A	0	alpha CAASFYKLSF	2	11
30	A	0	alpha CAASGDTDKLIF	1	12
31	A	0	alpha CAASGGAANKLIF	1	13

Showing 1 to 31 of 121,406 entries

Building a Dataframe `dat$normV2 <- dat$V2 * (100/sum(dat$V2))`



Building a Dataframe

dataframe



The screenshot shows an RStudio session with a code editor tab titled "TTAnalysis.Rmd" and a data viewer tab titled "dat". The data viewer displays a data frame with the following structure:

	ID	month	chain	V1	V2	CDR3Len	normV2
1	A	0	alpha	CAADQTGANNLFF	1	13	0.01476015
2	A	0	alpha	CAADSGGYQKVTF	1	13	0.01476015
3	A	0	alpha	CAADSGYSTLTF	27	12	0.39852399
4	A	0	alpha	CAADTGRRALTF	6	12	0.08856089
5	A	0	alpha	CAAFSGGYNKLIF	2	13	0.02952030
6	A	0	alpha	CAAIRGAGSYQLTF	1	14	0.01476015
7	A	0	alpha	CAAIRQAGTALIF	1	13	0.01476015
8	A	0	alpha	CAAIVAGTALIF	1	12	0.01476015
9	A	0	alpha	CAAIVIPGSQGNLIF	1	15	0.01476015
10	A	0	alpha	CAALWGYSSASKIIF	1	15	0.01476015
11	A	0	alpha	CAAMDSNMDNSNYQLIW	1	16	0.01476015
12	A	0	alpha	CAAPHRCGNEKLT	11	14	0.16236162
13	A	0	alpha	CAAPHRFGNENKLT	209	14	3.08487085
14	A	0	alpha	CAAPHRVGNEKLT	1	14	0.01476015
15	A	0	alpha	CAAPPRFGNEKLT	1	14	0.01476015
16	A	0	alpha	CAAPRGYSTLTF	1	12	0.01476015
17	A	0	alpha	CAAPTPLYNTDKLIF	2	15	0.02952030
18	A	0	alpha	CAARETGASKLTF	1	14	0.01476015
19	A	0	alpha	CAARETSGSRLTF	1	13	0.01476015
20	A	0	alpha	CAARKAETRGSRLTF	1	15	0.01476015
21	A	0	alpha	CAARKAETSGSRLTF	103	15	1.52029520
22	A	0	alpha	CAARKTGSRLTF	1	13	0.01476015
23	A	0	alpha	CAARRAETSGSRLTF	1	15	0.01476015
24	A	0	alpha	CAARVRDMRF	12	10	0.17712177
25	A	0	alpha	CAASAGGKLIF	1	11	0.01476015
26	A	0	alpha	CAASALPGNTPLVF	13	14	0.19188192
27	A	0	alpha	CAASDGNKLIF	1	11	0.01476015
28	A	0	alpha	CAASEERGATNKLIF	1	15	0.01476015
29	A	0	alpha	CAASFYKLSF	2	11	0.02952030
30	A	0	alpha	CAASGDTDKLIF	1	12	0.01476015
31	A	0	alpha	CAASGGAANKLIF	1	13	0.01476015

Showing 1 to 31 of 121,406 entries

Building a Dataframe

TTAnalysis.Rmd x dat x

Filter

	ID	month	chain	V1	V2	CDR3Len	normV2
1	A	0	alpha	CAADQTGANNLFF	1	13	0.01476015
2	A	0	alpha	CAADSGGYQKVTF	1	13	0.01476015
3	A	0	alpha	CAADSGYSTLTF	27	12	0.39852399
4	A	0	alpha	CAADTGRRAALTF	6	12	0.08856089
5	A	0	alpha	CAAFSGGYNKLF	2	13	0.02952030
6	A	0	alpha	CAAIRGAGSYQLTF	1	14	0.01476015
7	A	0	alpha	CAAIRQAGTALIF	1	13	0.01476015
8	A	0	alpha	CAAIVAGTALIF	1	12	0.01476015
9	A	0	alpha	CAAIVIPGSQGNLIF	1	15	0.01476015
10	A	0	alpha	CAALWGYSSASKIIF	1	15	0.01476015
11	A	0	alpha	CAAMDSNMDNSNYQLIW	1	16	0.01476015
12	A	0	alpha	CAAPHRCGNEKLTF	11	14	0.16236162
13	A	0	alpha	CAAPHRGFGNEKLTF	209	14	3.08487085
14	A	0	alpha	CAAPHRVGNGNEKLTF	1	14	0.01476015
15	A	0	alpha	CAAPPRFGNGNEKLTF	1	14	0.01476015
16	A	0	alpha	CAAPRGYSTLTF	1	12	0.01476015
17	A	0	alpha	CAAPTPLYNTDKLIF	2	15	0.02952030
18	A	0	alpha	CAARETGASKLTF	1	14	0.01476015
19	A	0	alpha	CAARETSGSRLTF	1	13	0.01476015
20	A	0	alpha	CAARKAETRGSRLTF	1	15	0.01476015
21	A	0	alpha	CAARKAETSGSRLTF	103	15	1.52029520
22	A	0	alpha	CAARKTSGSRLTF	1	13	0.01476015
23	A	0	alpha	CAARRAETSGSRLTF	1	15	0.01476015
24	A	0	alpha	CAARVRDMRIF	12	10	0.17712177
25	A	0	alpha	CAASAGGKLF	1	11	0.01476015
26	A	0	alpha	CAASALPGNTPLVF	13	14	0.19188192
27	A	0	alpha	CAASDGNKLVF	1	11	0.01476015
28	A	0	alpha	CAASEERGATNKLIF	1	15	0.01476015
29	A	0	alpha	CAASFQYKLSF	2	11	0.02952030
30	A	0	alpha	CAASGDTDKLIF	1	12	0.01476015
31	A	0	alpha	CAASGGAANKLIF	1	13	0.01476015

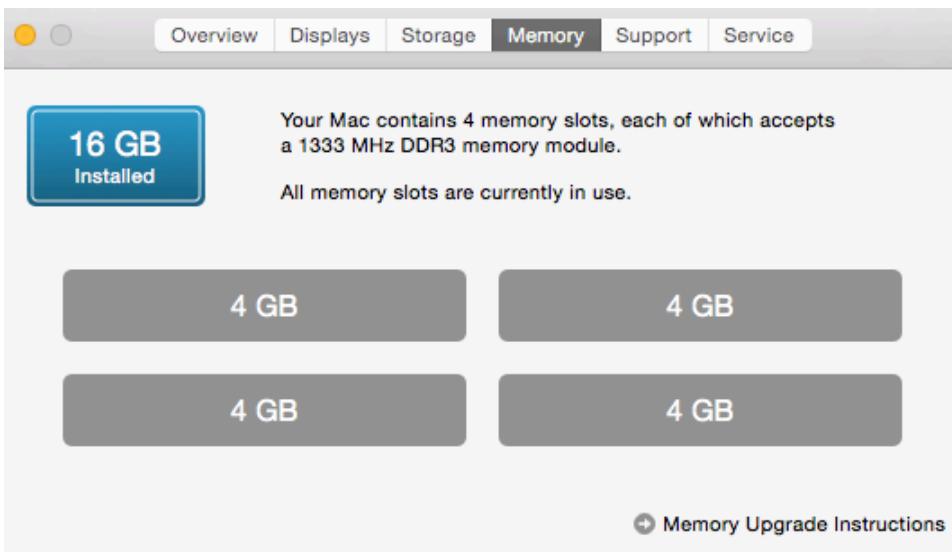
Showing 1 to 31 of 121,406 entries

+ Your analysis goes here!

Why this format?

Showing 1 to 31 of 121,406 entries

- Your data is loaded into RAM (fast memory)
- Your computer has 4 – 32 GB of available RAM



The screenshot shows the 'Memory' tab in Mac System Preferences. A blue box highlights '16 GB Installed'. Below it, text states: 'Your Mac contains 4 memory slots, each of which accepts a 1333 MHz DDR3 memory module.' It also says 'All memory slots are currently in use.' Four grey buttons below show '4 GB' for each slot.

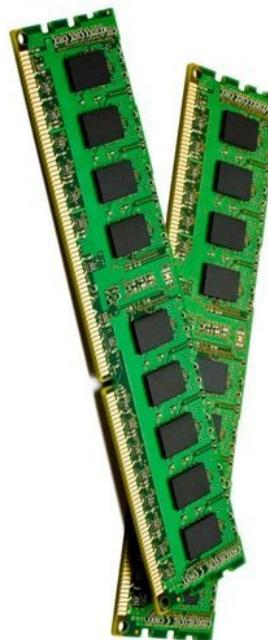
16 GB
Installed

Your Mac contains 4 memory slots, each of which accepts a 1333 MHz DDR3 memory module.
All memory slots are currently in use.

4 GB 4 GB

4 GB 4 GB

Memory Upgrade Instructions



Why this format?

- Every **character** in the dataframe **takes up some memory**
- We need to **condense** the **information** as much as possible to hold and analyse this much data. If something takes 1 hour to run. 10 * as much data will take 10 hours.

In a .csv file:

CAADQTGANNLFF, 250	17 bytes	58,823,529 copies
CAADQTGANNLFF, CAADQTGANNLFF, CAADQTGANNLFF, .	4.25 kilobytes	235,294 copies
.		
CAADQTGANNLFF		

Automated Analysis

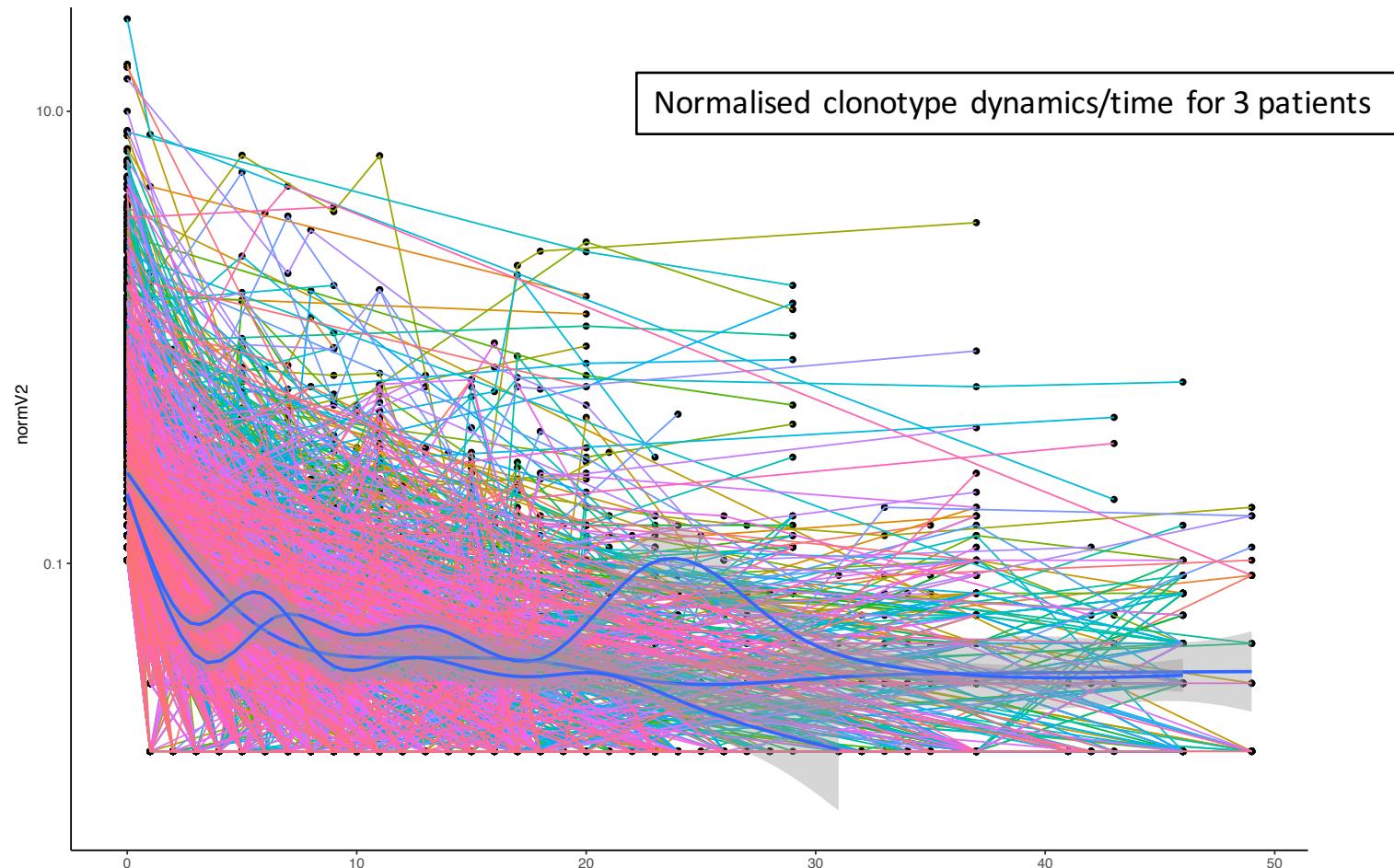
Q. How do we know the pipeline works?

A. The data are consistent with our biological understanding.

Automated Analysis

Q. How do we know the pipeline works?

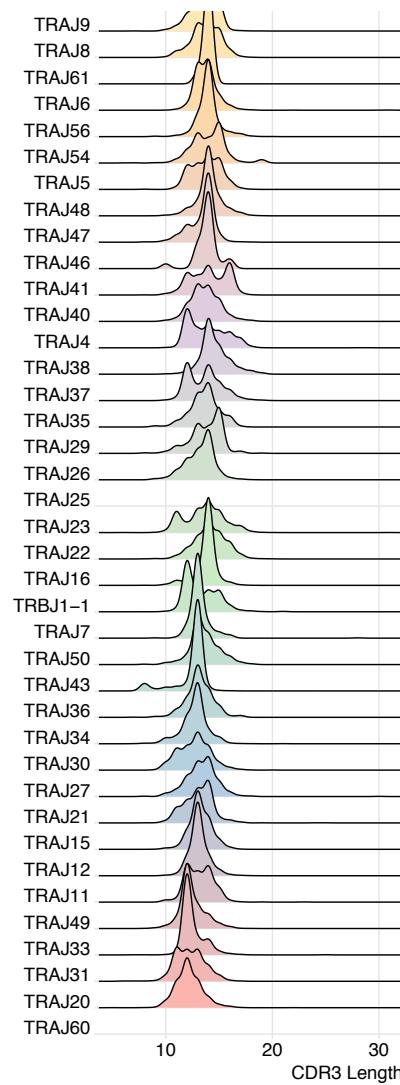
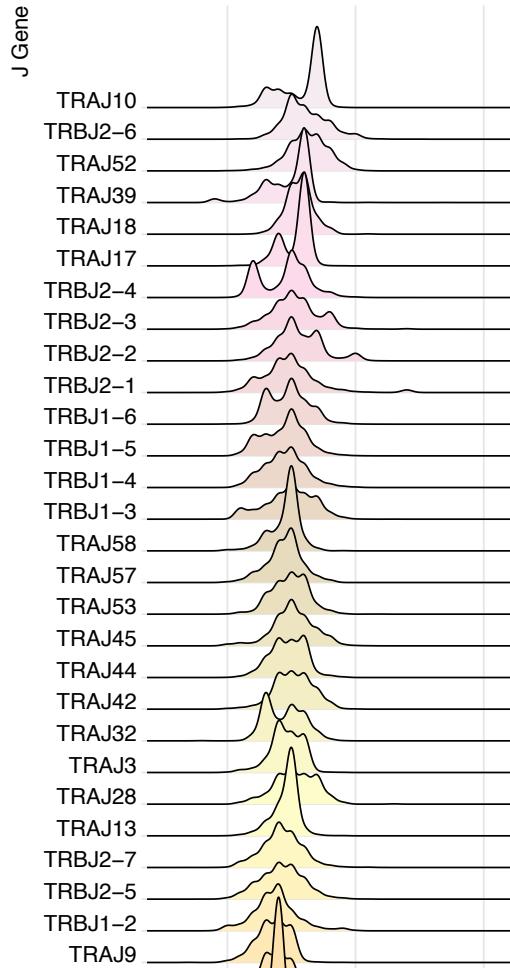
A. The data are consistent with our biological understanding.



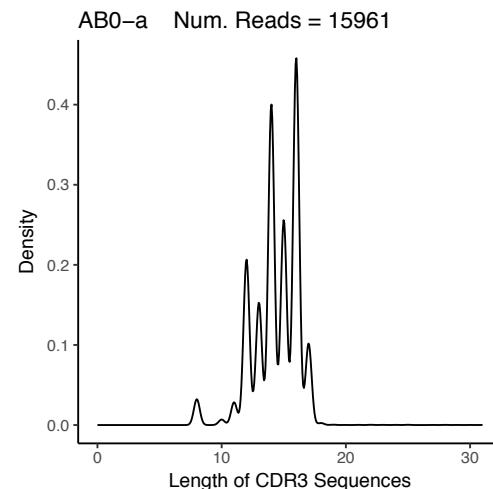
Automated Analysis



Using the data structure presented previously, we have automated:

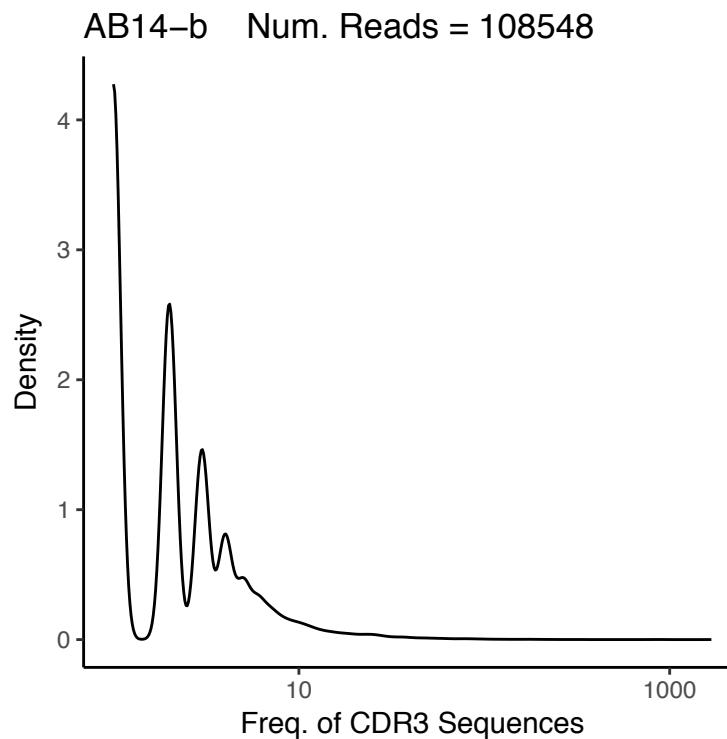
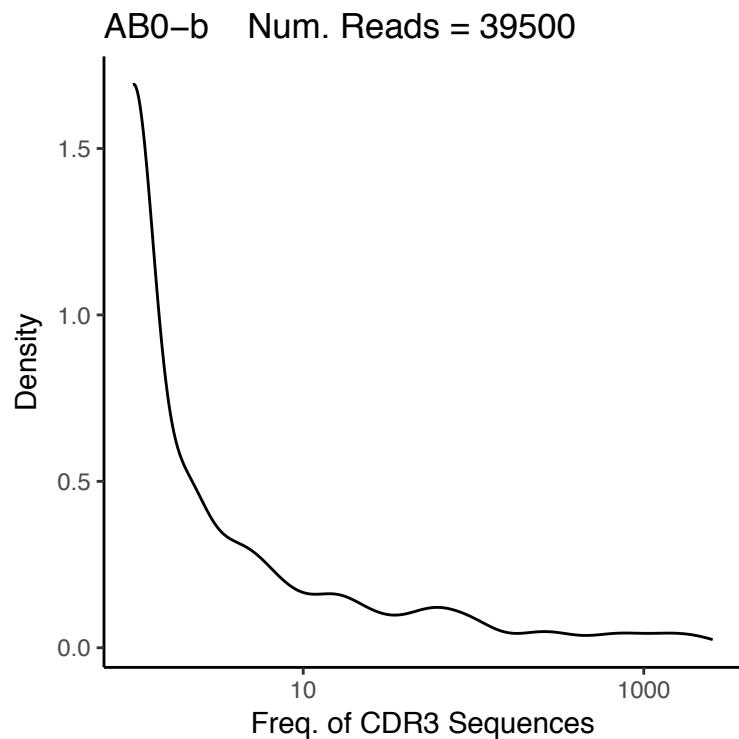


Spectratyping

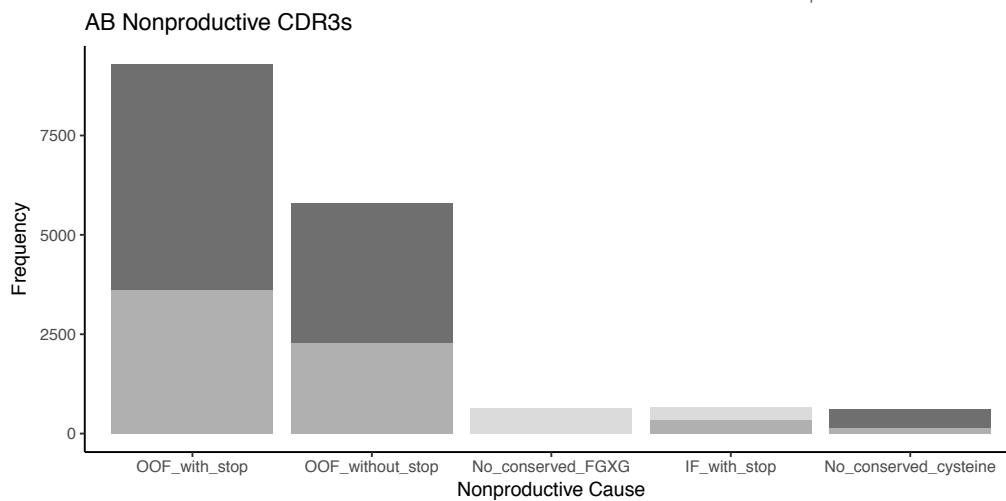
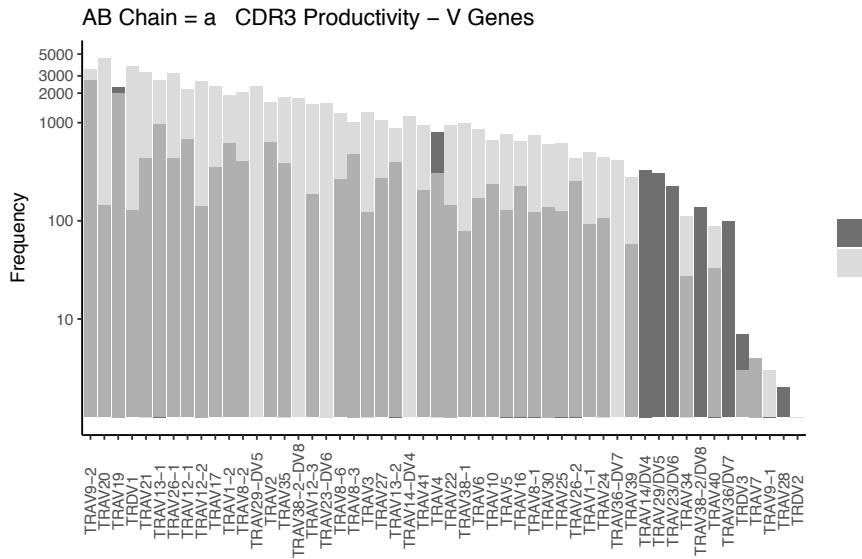
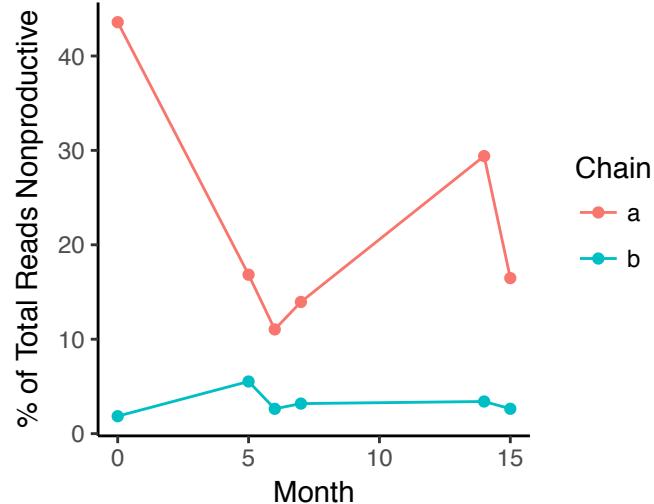


Using the data structure presented previously, we have automated:

Clonotype Frequency Density



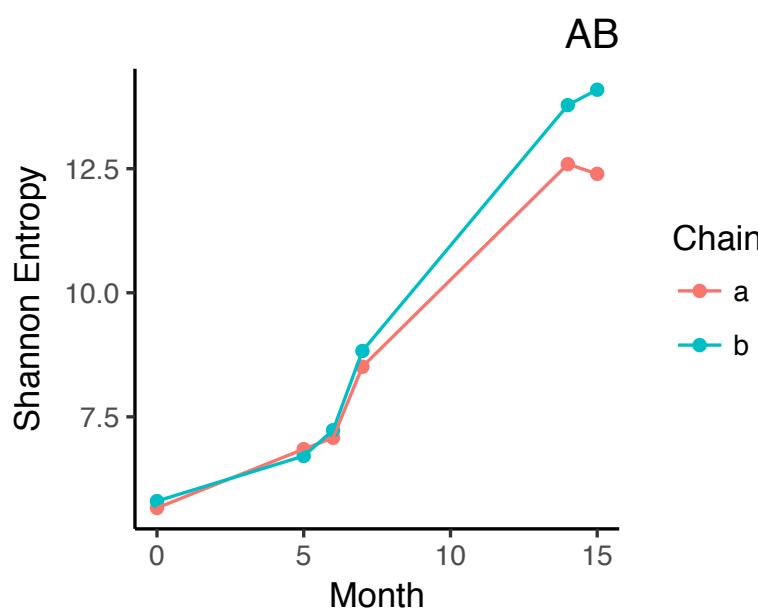
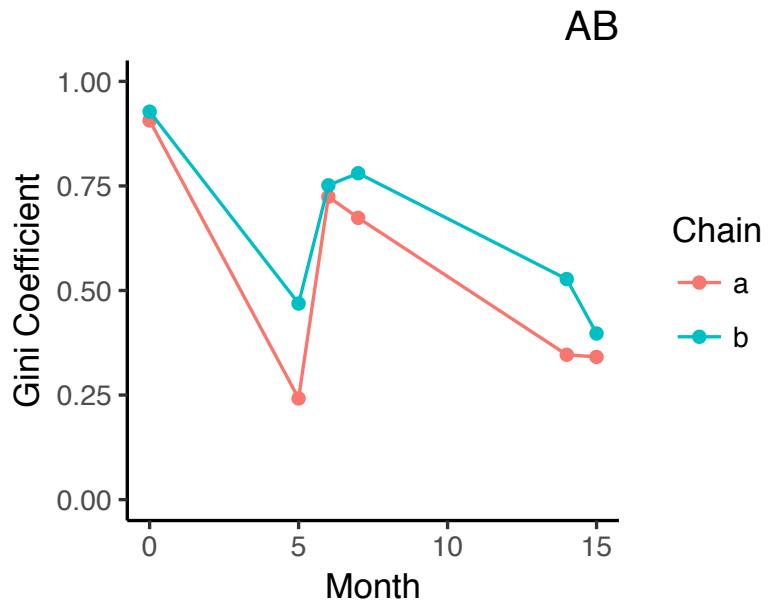
Using the data structure presented previously, we have automated:



Nonproductive analysis

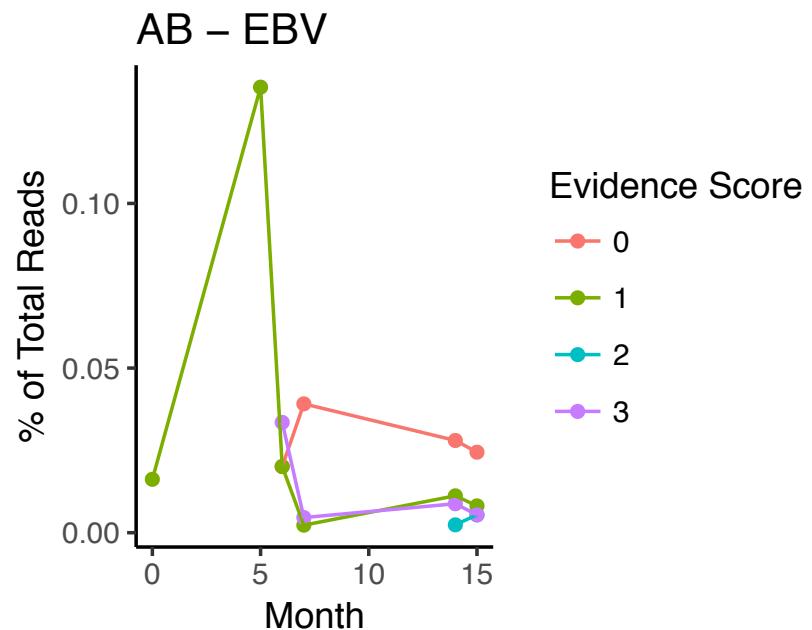
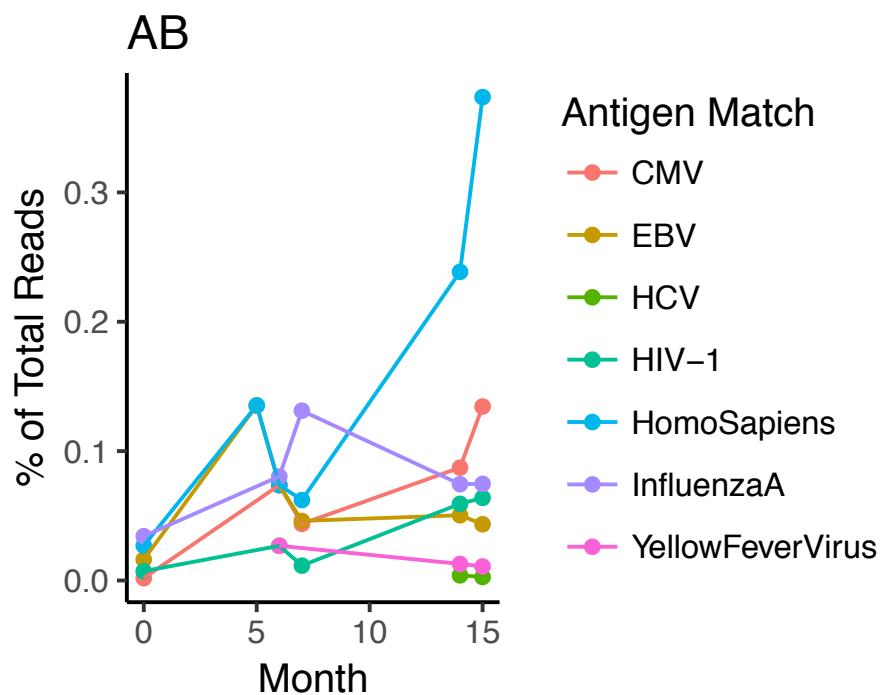
Using the data structure presented previously, we have automated:

Diversity Analysis



Using the data structure presented previously, we have automated:

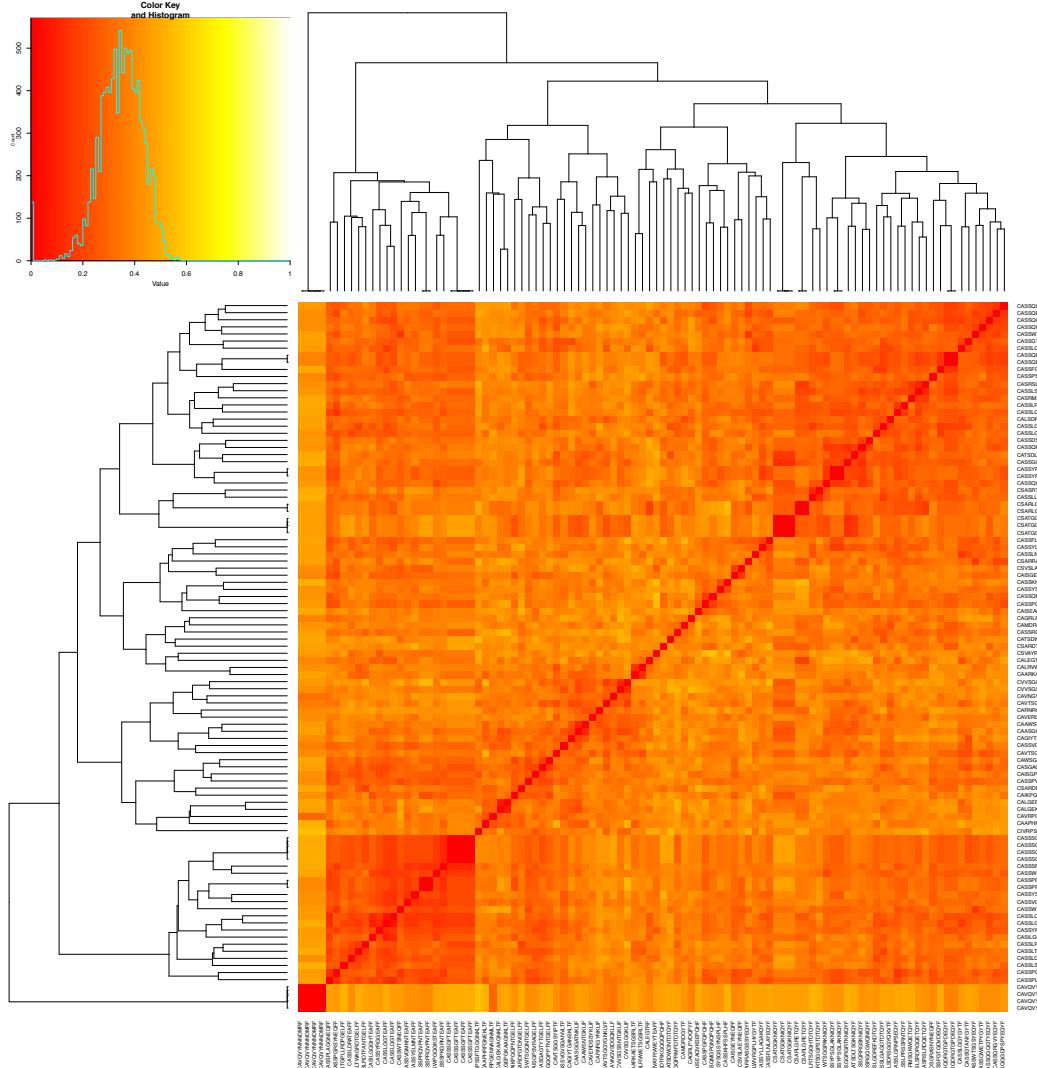
Antigen Specificity Database Matching
 (Although we don't recommend it!)



Automated Analysis

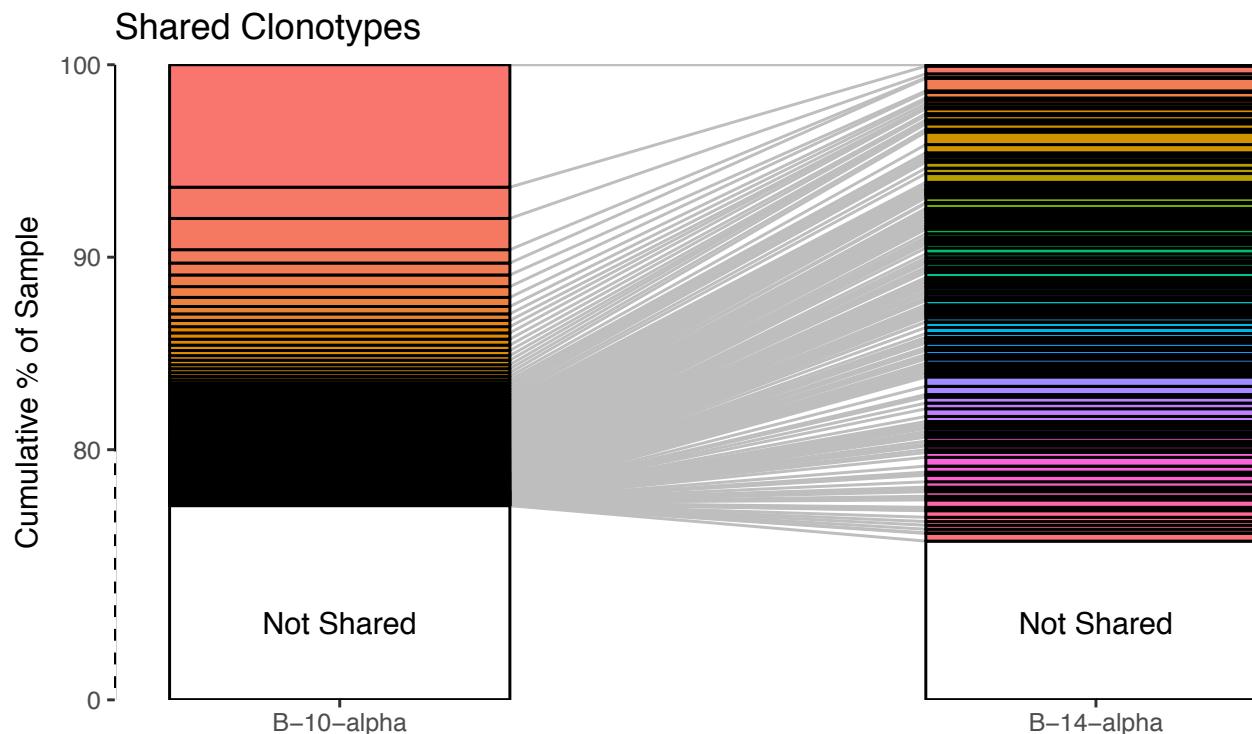
Using the data structure presented previously, we have automated:

Clonotype Similarity Analyses

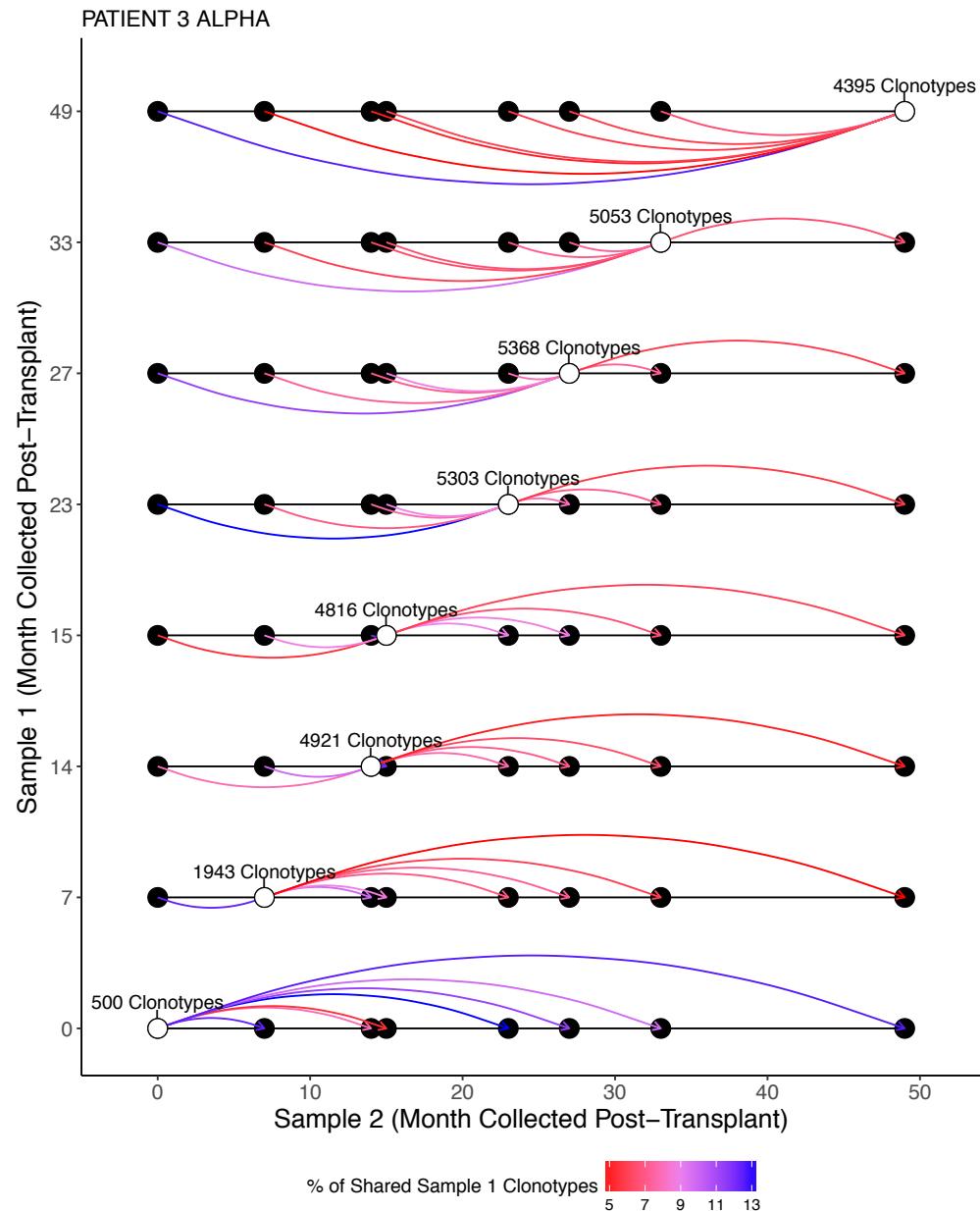


Making Your Own Analysis

- All the **code** is **open source** and available in the script.
- Custom analysis is highly encouraged.
- Draw the plot with pen and paper, break it down into manageable steps to program.



Automated Analysis



... It is easier then you think:

32 lines of code on-top of dataframe

4 lines to automate

~ $\frac{3}{4}$ of a days work

Common Issues

The errors produced by R are not exactly intuitive

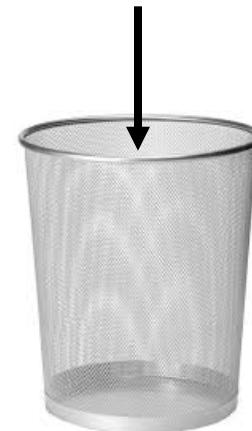
```
> #Checks the R1 file against file names in the info file
> fileName <- strsplit(fastQFileR1,'/')[[1]]
Error in strsplit(fastQFileR1, "/")[[1]] : subscript out of bounds
> |
```

If you are seeing errors, then it is likely caused by the following issues:

- 1). You've given it the wrong file path (R can be very picky).
- 2). You have no (or very few) reads in your data (check the CDR3 translator output file sizes).
- 3). AutoTCR or one of its dependencies has not been installed correctly.

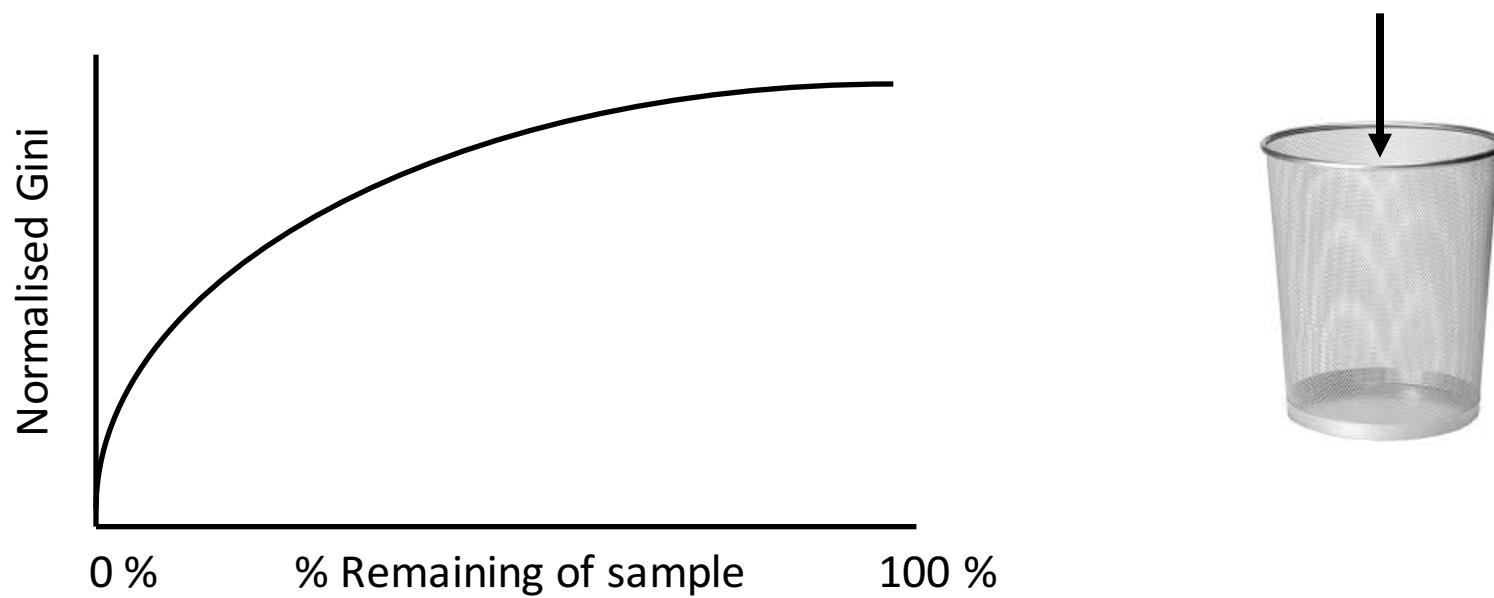
A Quick Note About Subsampling

- Randomly subsampling all of your samples to the **same depth** is **essential** for comparison.
- However, **subsampling too far down** (i.e. < 10 %) is **just as bad** as not subsampling your data. You are literally **throwing away** > 90 % of your data!



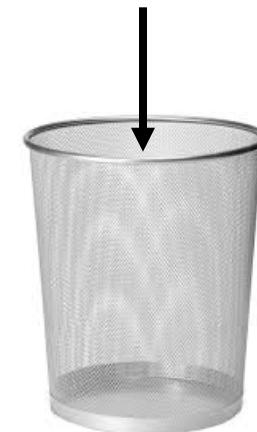
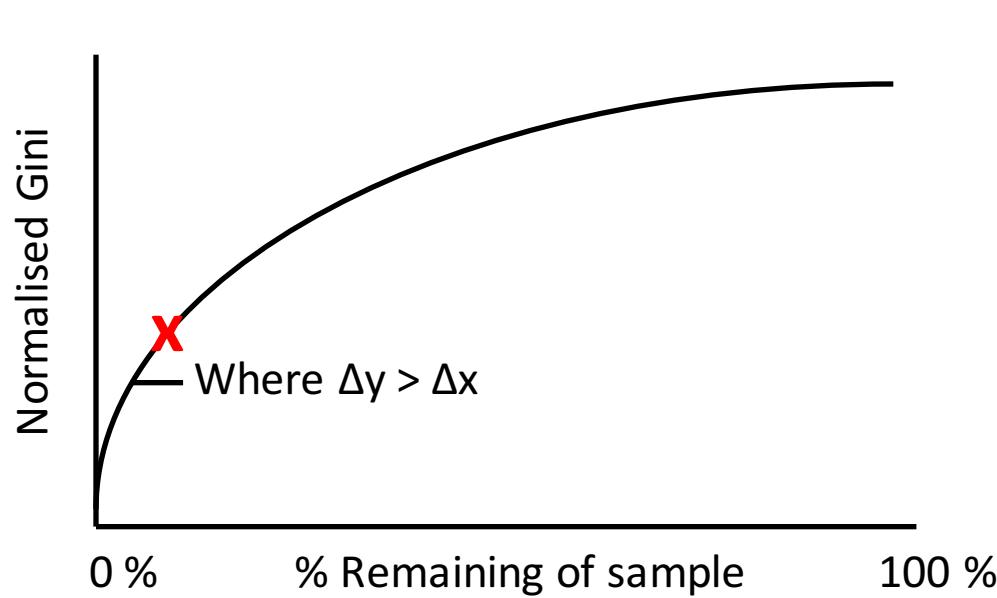
A Quick Note About Subsampling

- Randomly subsampling all of your samples to the **same depth** is **essential** for comparison.
- However, **subsampling too far down** (i.e. < 10 %) is **just as bad** as not subsampling your data. You are literally **throwing away** > 90 % of your data!



A Quick Note About Subsampling

- Randomly subsampling all of your samples to the **same depth** is **essential** for comparison.
- However, **subsampling too far down** (i.e. < 10 %) is **just as bad** as not subsampling your data. You are literally **throwing away** > 90 % of your data!



A Quick Note About Subsampling

X (position) is very variable between each sample

We calculate the value of **X** for each sample

Next we calculate a subsampling depth that includes the maximum number of samples

A Quick Note About Subsampling

X (position) is very variable between each sample

We calculate the value of **X** for each sample

Next we calculate a subsampling depth that includes the maximum number of samples

If this isn't considered, you risk artificially separating patient populations



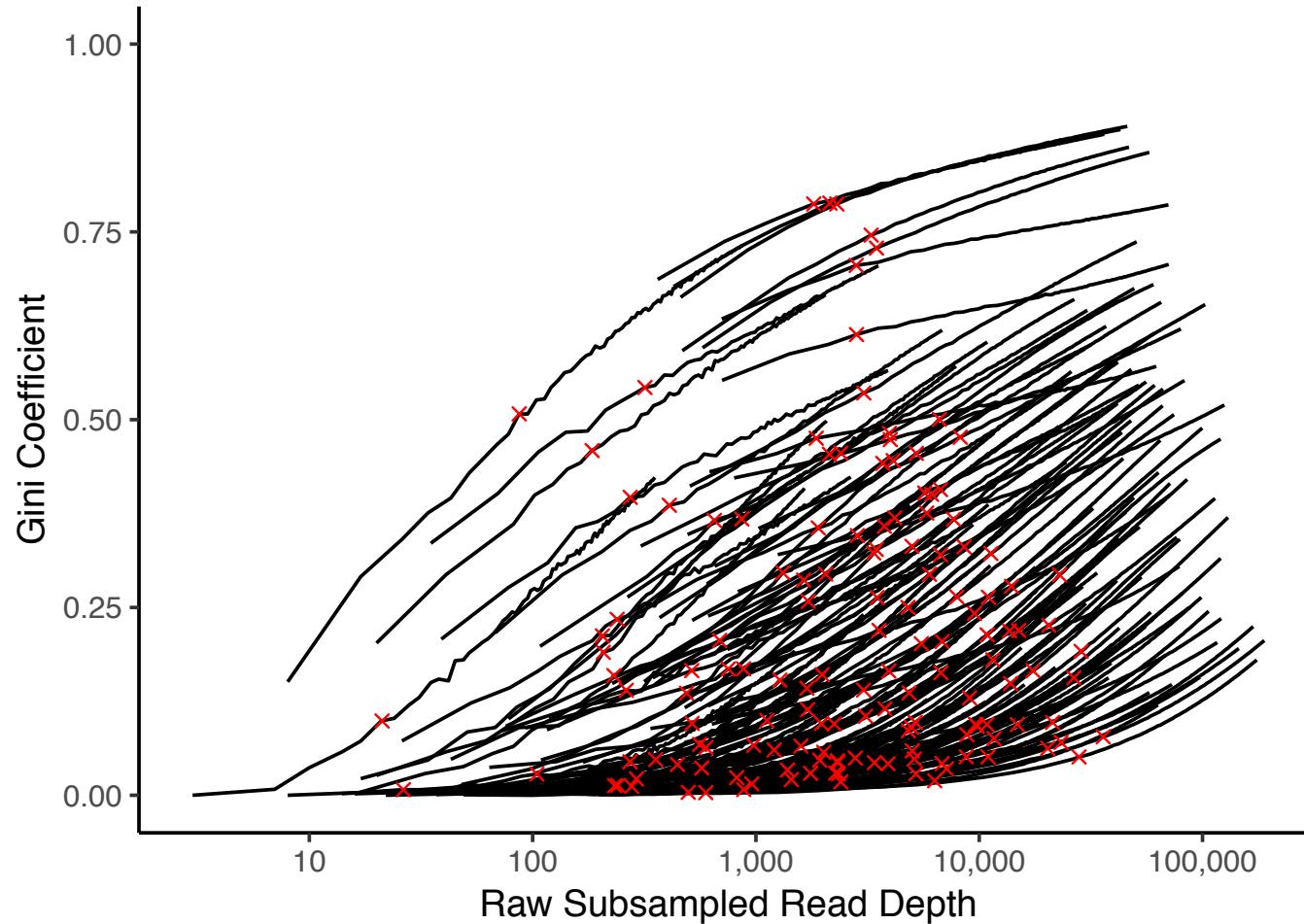
Control

- Normal cell count
- Large number of reads
- Already good diversity
- Needs sig. subsampling to be compared with...

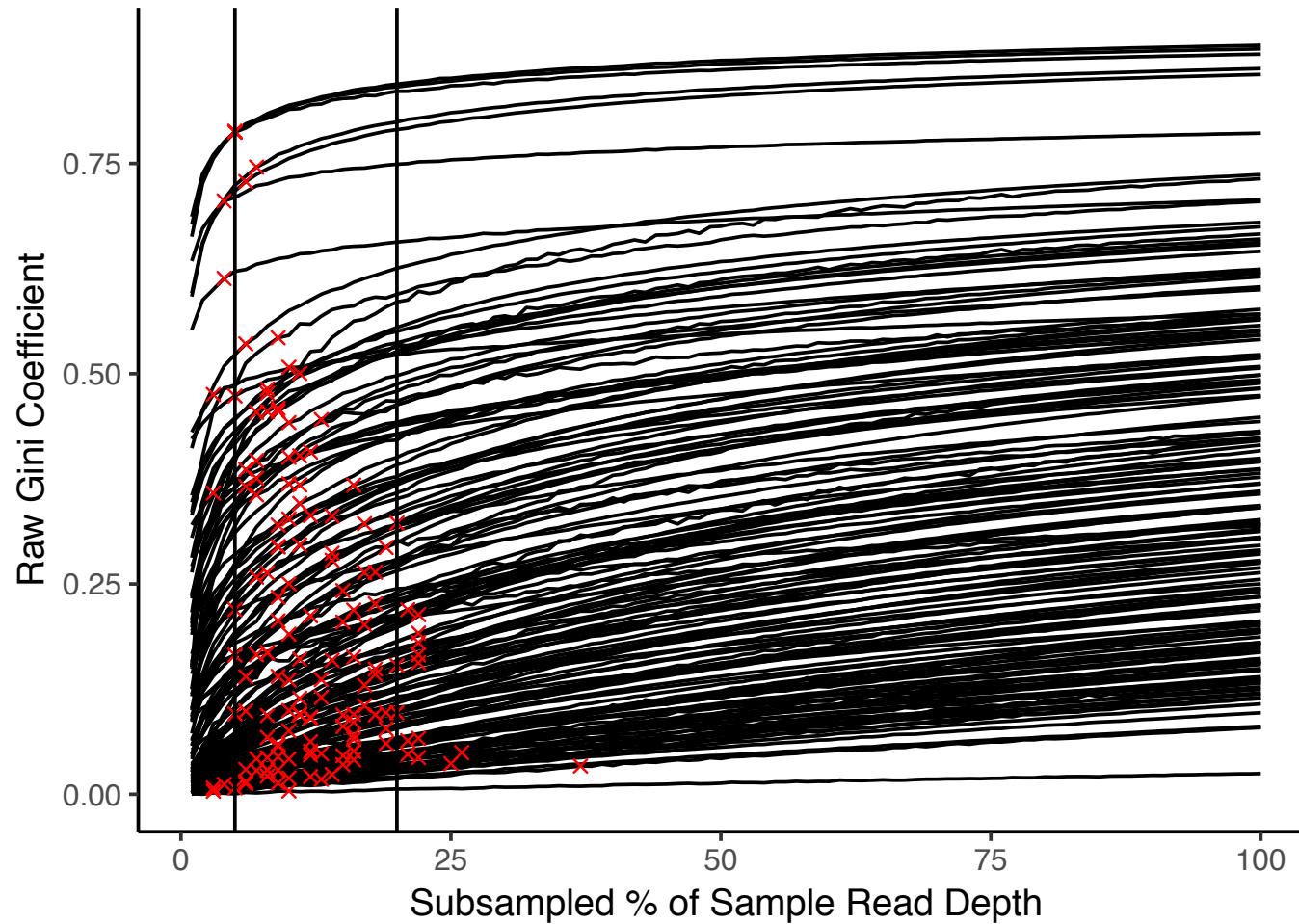


Patient

A Quick Note About Subsampling



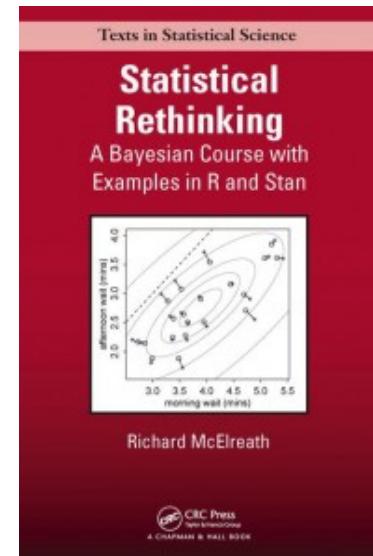
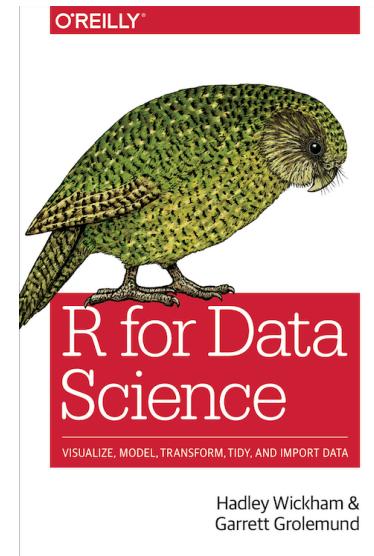
A Quick Note About Subsampling



Normalising is key!

Resources

- Both of us feel strongly that **equal time** between lab and analysis is essential.
- <http://r4ds.had.co.nz> (Hadley Wickham)
- Statistical Rethinking (Richard McElreath)
- <https://www.codecademy.com>
- <https://learnpythonthehardway.org>
- <https://www.kaggle.com>
- <https://www.tidyverse.org>
- <https://pandas.pydata.org/>



(Last tip – backup everything)

Suggested Future Work

(The things we ran out of time to do)

- We worked with the data and resources available.
- Thinking of the next steps is key (lab techniques + data analysis)

Lab

- Automated library prep (Hamilton)
- B cell sequencing
- High-throughput sequencing (MiSeq -> NextSeq/HiSeq?)
- Single cell sequencing

Data Analysis

- Further streamlining of AutoTCR/Package publishing
- Interactive analysis
- Integrating diversity w/ cell counts
- Machine learning
- Mixture models

Acknowledgements

Athina Soragia Gkazi (Lab)
Stuart Adams (Lab)
Katrine Schou-Sandgaard (Lab)
Rose Gkouleli (Lab)
Lana Mhaldien (Lab)
Theres Oakes (Lab)
Eleanor Watt (Lab)
Robin Callard (PI)
Graham Davies (PI)
Joseph Standing (Supervisor)
Nigel Klein (Supervisor)
Judith Breuer (Supervisor)
Jamie Heather (Software)
Richard Goldstein (Modelling)
Watjana Lilaonitkul (Modelling)
Theresa Attenborough (Data)
John Booth (Data)
Benny Chain (Data)

Questions?

