Adaptive Immune Receptor Repertoire Sequencing (AIRRseq)

Barbera van Schaik Bioinformatics Laboratory Epidemiology & Data Science Amsterdam UMC b.d.vanschaik@amsterdamumc.nl

Overview

Introduction AIRRseq

- * T- and B-cell receptors
- * VDJ recombination
- * Clonal expansion and somatic hypermutation
- * Class switching
- * Repertoire sequencing

Data analysis

* Pre-processing

- * Identify V, (D), J and C regions
- * Identify CDR3 region
- * Determine clones and their frequency

Follow-up analysis

- * Dominant clones
- * Compare samples
- * Diversity of samples
- * Shared clones between samples

Computer systems for data analysis

- * Clusters, cloud, linux
- * Overview of computer lab

Innate and adaptive immune system



Nature Reviews | Cancer

T and B cells

- 1. Part of the adaptive immune system
- 2. Slower, but **specific** immune response
- 3. Trick is the high variety of Tand B-cells: immune receptors



Immune receptors: V, D, J recombination

- BCR: heavy (heavy) and light chains (kappa, lambda)
- TCR: beta/alpha and gamma/delta chains



⁵ https://digitalworldbiology.com/blog/immunoprofiling-how-it-works

Human Antigen Receptor Loci



[©] Digital World Biology, LLC

https://digitalworldbiology.com/blog/immunoprofiling-how-it-works



https://digitalworldbiology.com/blog/immunoprofiling-how-it-works



Ralph et al (2016), PloS Computational Biology

Clonal expansion and somatic hypermutation of B cells



Wikipedia

Class switching

Genes in heavy chain locus of an IgM expressing B cell



Repertoire sequencing

- DNA
- RNA
- Bulk (one chain)
- Single cell

One of the protocols

STEP 1: Specific reverse transcription with UMI-tagged primers



STEP 2: Exonuclease I treatment



STEP 3: Multiplex PCR amplification (35 cycles) + 2x AMpure XP clean up



STEP 4: Indexing PCR (5 cycles) + 2x AMpure XP clean up



Pollastro et al (2021), Front Immunol.

OUTPUT: Amplicon ready for sequencing on Illumina MiSeq



Unique Molecular Identifiers (UMIs)



Unique Molecular Identifiers (UMIs)



https://www.ecseq.com/support/ngs/how-can-unique-molecular-identifiers-help-to-reduce-quantitative-biases



Data analysis

Pre-processing

- Identify barcodes
- Identify UMIs
- Mask primer sequences



Identify V, (D), J and C regions

$ \begin{array}{c} \mathbb{P} \mathbb{P} \mathbb{Q} \mathbb{P} \mathbb{Q} \mathbb{L} \mathbb{P} \mathbb{Q} \mathbb{Q} \mathbb{P} \mathbb{Q} \mathbb{Q} \mathbb{P} \mathbb{Q} \mathbb$	L R L S C A A S G F T F S CTGAGACTCTCCTGTGCAGCCTCTGGATTCACCTTTAGC 213 90 L R L S C A A S G F T F S 90
Query_1 124 GAGGTGCAGCTGTTGGAGGTCTGGGGGGAGGCTTGGTACAGCCTGGGGGGGG	L R L S C A A S G F T F S
V 93.2% (275/295) <u>IGHV3-23*01</u> 1 E V Q L L E S G G G L V Q P G G S V 93.2% (275/295) <u>IGHV3-23D*01</u> 1	LRLSCAASGFTFS 90
EVQLLESGGGLVQPGGS V 93.2% (275/295) TCHV3-23D*01 1	LRLSCAASGFTFS
V 93.2% (275/295) TCHV3-23D*01 1	
V 92.9% (274/295) IGHV3-23D*02 1G	
GTFR2-IMGTFR2-IMGT	> <cdr2-imgt><</cdr2-imgt>
TYAMSWVRQAPGKGLEW	V S T I T G D D G T T Y Y
Query 1 214 ACCTATGCCATGAGCTGGGTCCGCCAGGCTCCAGGAAAGGGGCTGGAGTG	GTCTCAACTATTACTGGTGATGATGGTACCACATACTAC 303
V 93.28 (275/295) IGHV3-23*01 91 .G	GGAGGG
SYAMSWV BOAPGKGLEW	V S A I S G S G G S T Y Y
V 93.28 (275/295) TGHV3-23D*01 91 .G.	
V 92.98 (274/295) TGHV3-23D+02.91 G	G G AG G G 180
	Internet in Addition internet in the second
2 G V	THOM
• • • • • • • • • • • • • • • • • • •	TUNE ON NOT BARD
	ACAGEGEAMORGCAAABGEACACCCCCAACAC 262
U 03 25 (27E/28E) TOTY 23401 101 C 8	ACADIGINICIOCANNIGANCOULTICAGNOCCONNUM. 373
A D S V K G R F T I S R D N S K N	TLYLOMNSLRARD
V 93.28 (275/295) IGHV3-23D*01 181 G.A.	
V 92.98 (274/295) TGRV3-23D*02 181 G.A	
CDR3-IMGT	>
A A V Y Y C A K G R C G D S W C S	G F D C W G Q G I L V T V
Query 1 394 GCGGCCGTGTATTACTGTGCGAAAGGGCGTTGTGGTGATAGCTGGTGCTC	GGCTTTGACTGCTGGGGGCCAGGGAATCCTGGTCACCGTC 483
V 93.2% (275/295) IGHV3-23*01 271 AAA	295
TAVYYCAR	
V 93.2% (275/295) IGHV3-23D*01 271 AA	295
V 92.9% (274/295) IGHV3-23D*02 271 AA.	
D 100.0% (8/8) IGHD2-21*01 12	
D 100.0% (8/8) IGHD2-21*01 7	14
D 100_0% (8/8) TGHD2-21*02 7	14
J 95.5% (42/44) TGHJ4*02 5	

					3 5			
			Query_1	484	TCCTCAG	490		
J	95.5%	(42/44)	IGHJ4*02	42		48		
J	93.2%	(41/44)	IGHJ4*01	42		48		
J	97.18	(34/35)	IGHJ5*02	45		51		

. .

*Toby et al (2016), BMC Bioinformatics*¹⁸

Identify CDR3 sequence

- Two common methods for this
 - Deduce from alignment
 - Search for conserved motifs at end of V and start of J region

			NDN	J		
FW1	CDR1	FW2	CDR2	FW 3	CDR3	FW4

Define (sub)clones

Sub clones

V, J and CDR3-peptides

V, J and CDR3-nucleotides

Entire nucleotide sequence

Entire protein sequence

Clonal families (clones)

1 amino acid difference

85% sequence identity

Dynamic threshold



https://changeo.readthedocs.io/en/stable/examples/cloning.html

20

Follow-up analysis

Dominant clones



- Defined as clones with percentage reads (or UMIs)
 > 0.5% compared to total sample
- This threshold was supported by naive sorted T cells

Compare samples

в







Α

Pollastro et al (2019), ARD

Diversity

Diversity



Many different plants + Equal frequencies = diverse



Few types of plants + Dominated by lawn = less diverse

Diversity indices

Shannon entropy

$$H_1 = -\sum_s p(s) \ln p(s).$$

Renyi entropy

$$H_{\beta} = \frac{1}{1-\beta} \ln \left[\sum_{s} p(s)^{\beta} \right]$$

Hill diversity

$$D_{\beta} = \exp[H_{\beta}].$$

Inverse Simpson

$$D_2 = 1 / \sum_s p(s)^2.$$

Gini-Simpson

 $1 - 1/D_2$.

p(s) is probability, frequency or abundance of species

D0 is species richness

²⁶ Mora and Walczak, 2016

Example immune cells



Generalized formula for entropy and diversity

- $\beta = 0$: insensitive to frequency
- ß = low: emphasis on rare species
- ß = high: emphasis on dominant species

$$H_{\beta} = \frac{1}{1-\beta} \ln \left[\sum_{s} p(s)^{\beta} \right] \qquad \qquad D_{\beta} = \exp[H_{\beta}].$$

Renyi

Hill diversity

Diversity profiles



Greiff et al (2015), Genome Medicine

Shared clones

Similarity indices

- Inverse of distance measure
- E.g. determine overlap in species richness (species count, without abundance)
- Several similarity indices exist that do take abundance into account



Community 1 A: 25% B: 25% C: 25% D: 25%





Quality control: Contamination (similarity)

- Pairwise ² comparison between all samples
- 0=no overlap at all (dark blue)
- 1=identical samples (yellow)



Quality control: Contamination (similarity)

- Black squares: samples from same patient
- Similarity is expected within patient
- Similarity across patients indicate problems (contamination, labelling errors, etc)



Similarity indices

Jaccard index (occurrence of species) J(A,B) = intersection_AB / (union_AB - intersection_AB)

Sorensen index (occurrence of species) S(A,B) = 2 * intersection_AB / (set_A + set_B)

Bray-Curtis index (occurrence and frequency of species) $BC_{ij} = 2 * C_{ij} / (S_i + S_j)$

 $\mathbf{C}_{_{ij}}$ is sum of values for species in common between both sites

S_i is total number of species at site i

S_i is total number of species at site j



union

setA setB

union

Look at shared clones

Compare all samples pairwise Look for shared clones (same CDR3 sequences) in both samples Calculate impact of the shared clones in both samples: Impact = sum(frequency) / total frequency in sample

Assumption:

Contamination direction is from sample with highest impact to lowest impact

Shared clone:

Purple = same patient Blue = different patient

Impact of target sample



Impact of source sample

Shared clone

Only show different patient pair

Impact of target sample



Impact of source sample

Shared clones

File Edit View Insert Format Sheet Data Tools Window Help

- 7 -

Δ1

~1															
	A	В	С	D	E	F	G	Н		J	K	L	М	Ν	Ī
1		index	source	target	impact_source	impact_target	patient_source	patient_target	₄me patie	nt					
2	0	8600	S074-283_S63	S074-093_S104	91.8363838134388	89.5445790769681	53	31	different						
3	1	21516	S074-078_S103	S074-273_S62	87.1523661894775	84.7566716572258	53	43	different						
4	2	831	S074-290_S139	S074-279_S121	77.9119642794361	77.552679170138	22	22	same						Ī
5	3	677	S074-290_S139	S074-004_S18	79.3587174348698	75.132630038417	22	22	same						Ī
6	4	676	S074-279_S121	S074-004_S18	78.4731676589519	74.4008781023233	22	22	same						
7	5	654	S074-181_S148	S074-004_S18	74.0881900098438	74.0045124702726	22	22	same						
8	6	7084	S074-290_S139	S074-181_S148	77.9365749041944	73.1593427395946	22	22	same						Ī
9	7	7083	S074-279_S121	S074-181_S148	77.030459431451	72.8867462581085	22	22	same						Ī
10	8	21527	S074-209_S204	S074-167_S110	58.6370428116676	57.3129770992366	43	43	same						Ī
11	9	21517	S074-167_S110	S074-078_S103	58.8187022900763	52.0269900610924	43	53	different						
12	10	21518	S074-209_S204	S074-078_S103	60.4706684856753	51.8665086167594	43	53	different						Ī
13	11	17694	S074-166_S109	S074-188_S154	52.3903188828679	51.5830611512829	83	83	same						Ī
14	12	21514	S074-167_S110	S074-273_S62	59.2786259541985	51.3636574109803	43	43	same						Ī
15	13	21515	S074-209_S204	S074-273_S62	61.1365555772104	51.1367670437883	43	43	same						Ī
16	14	20078	S074-175_S129	S074-238_S35	48.9532162509751	47.7550297952682	51	51	same						Ī
17	15	20777	S074-146_S67	S074-175_S129	51.2577526098848	47.5954544496216	51	51	same						
18	16	20077	S074-146_S67	S074-238_S35	51.6333448900189	46.8228214053927	51	51	same						Ī
19	17	17870	S074-188_S154	S074-303_S182	51.8832480601129	46.5720363980907	83	83	same						Î
20	18	17960	CO74 166 C100	CO74 202 C192	52 6026460262224	16 2021512226001	02	02	camo						Î

Computer clusters, cloud, linux

Handling sequencing data





Small cluster



Option: buy a computer cluster





Option: cloud computing



Often you do not need computing power all the time Compute by demand Compute resources are shared

Computer exercises

- Data analysis on the SurfSara cloud
- Analyze a public AIRRseq dataset
 - Pairwise assembly
 - Identify V, J and CDR3
 - Group related sequences
 - Count dominant clones



Overview

Introduction AIRRseq

- * T- and B-cell receptors
- * VDJ recombination
- * Clonal expansion and somatic hypermutation
- * Class switching
- * Repertoire sequencing

Data analysis

* Pre-processing

- * Identify V, (D), J and C regions
- * Identify CDR3 region
- * Determine clones and their frequency

Follow-up analysis

- * Dominant clones
- * Compare samples
- * Diversity of samples
- * Shared clones between samples

Computer systems for data analysis

- * Clusters, cloud, linux
- * Overview of computer lab