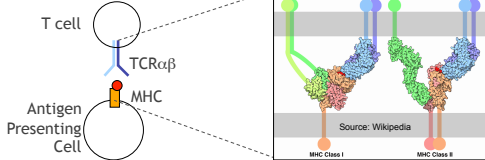


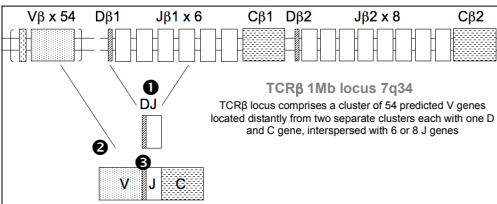
## T Cells

- T cells are central to cell-mediated immunity. They are distinguished from other lymphocytes, such as B cells, by the presence of the **T cell receptor (TCR)** on its cell surface
- T cells respond to antigenic peptides presented at the surface of nucleated cells by MHC



## T Cell Receptors (TCR)

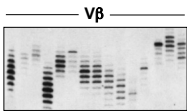
- TCR is a heterodimer and each chain includes a variable (V), joining (J) and a constant (C) segment. TCR $\beta$  also has a diversity (D) domain
- Structural TCR diversity necessary for recognition of enormous number of potential antigens generated by:
  - DJ gene recombination at the DNA level
  - One of the V genes joins DJ and the intermediary DNA is deleted
  - During the gene rearrangements, random base addition at the V-D-J junction and frequent base deletion V-3', 5'-J and either side of D yield the CDR3
- At V-(D)-J junction, the 3rd complementarity-determining region (CDR3) of TCR interacts directly with MHC-bound antigenic peptides
- CDR3 is the most critical TCR structure in epitope recognition
  - >10<sup>18</sup> theoretically possible  $\alpha\beta$ TCRs and est. ~10<sup>7</sup> T cell clonotypes in a given individual at a given time (Arestia et al. 1999)



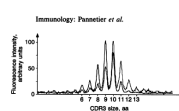
*Genomic re-arrangement creates a diverse T cell metagenome in every individual*

## Current TCR Profiling

- Spectratyping: Analysis of TCR  $\beta$ -chain repertoire complexity based on CDR3 length diversity within V $\beta$  gene families. PCR-based



Spectratype from adult blood lymphocytes (Gorski et al. 1994)

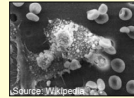


Immunoscope from mouse thymus (Pannetier et al. 1993)

- Spectratypes (immunoscopes) provide two levels of information: band number and band (fluorescence) intensity pattern
- Do not allow specific identification of a single T cell clone
- TCR sequence repertoire previously inaccessible due to cost of sequencing

## Objectives

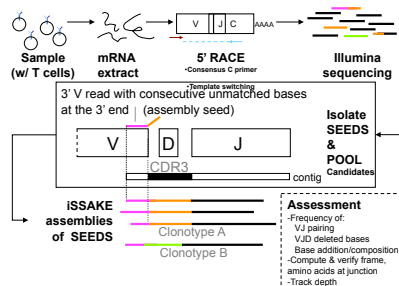
- Develop an approach for sequence profiling entire human T cell repertoire from blood
  - Deep illumina sequencing of T cell receptor variable region
  - De novo assembly using ISSAKE (Warren et al. 2007)



## Motivations

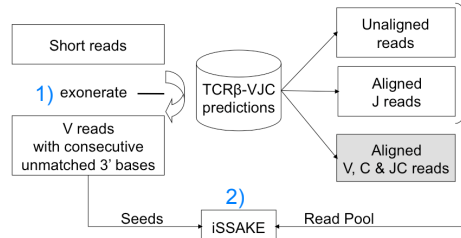
- Current profiling methods not quantitative
  - Current sequence profiling are low-scale and costly
    - hundreds clonotypes Vs. ~10<sup>7</sup> possible
    - 10-fold coverage of 1 Million 300 bp target sequences: \$3M Vs. \$3K (Sanger Vs. illumina)
- Applications include profiling cellular immune response to cancer, infection, vaccination, etc.
  - Compare profiles in sick Vs. healthy / recurrent Vs. non-recurrent tumors
  - Reconstitution of immune system after bone marrow replacement
- Next-generation sequencing make large-scale sequence profiling of T cell metagenomes a possibility, not without challenges due to read error and size

## TCR mRNA Sequencing



## TCR Assembly Strategy

- Designed a strategy for TCR sequence-profiling with short-read data
  - Method relies on: 1) sequence alignment 2) de novo assembly
- Reads aligning to the end of known TRBV genes and having consecutive unmatched bases in the adjacent CDR3 are used to seed ISSAKE de novo assemblies of non-templated CDR3



## Modeling

- Models of sequence diversity for the TCR beta-chain CDR3 region were built using empirical data and used to simulate, at random, distinct TCR clonotypes at 1-20 parts per million

## Simulation

- Using simulated TCR $\beta$  (sTCR $\beta$ ) sequences, we randomly created 20 million 36nt reads having 1%-2% random error, 20 million 42nt or 50nt reads having 1% random error and 20 million 36nt reads with 1% error modeled on real short read data

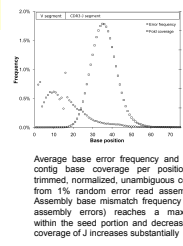
## Results

- With assembled 36nt reads we detect over 51% and 63% of rare (1 p.p.m.) clonotypes using a random or modeled error distribution

Over 99% of the more abundant clonotypes (6 p.p.m. or higher) can be detected using either error distribution

### TCR $\beta$ reconstructions - general stats

Reads	Error (%)	Mean reads	Short contigs (42 and 50nt) experiments	Long contigs (36nt) experiments	Classifications, but not optimal	A	B
36	1.02	1,399,140	498,583	36,352	1,064,399	1,098	2,847
36	1.0	1,755,437	205,618	11,230	1,478,299	157	89
36	1.3	1,714,470	339,423	16,400	1,342,285	238	927
36	2.0	1,678,905	441,317	22,618	1,215,270	302	1,198
42	1.0	2,448,325	380,847	21,056	2,046,549	239	818
50	1.0	3,114,789	361,076	34,528	2,428,086	174	459
50	1.0	3,114,789	361,076	34,528	2,428,086	174	459



Sensitivity, accuracy and clonotype frequency estimates

Error Distribution:						Read Length:					
Number of TCR $\beta$ -CDR3 characterized by ISSAKE contigs <sup>a</sup> (sequencing)			Accuracy (%)			Number of TCR $\beta$ -CDR3 characterized by ISSAKE contigs <sup>a</sup> (sequencing)			Accuracy (%)		
ppm	1.0	1.0	1.0	1.0	1.0	ppm	1.0	1.0	1.0	1.0	1.0
random	modeled	random	modeled	random	modeled	random	modeled	random	modeled	random	modeled
1	70,240	36,564	99.68	99.01	2.0	1.8	1	90,210	104,155	99.74	99.75
2	9,111	8,100	99.96	99.34	3.7	2.5	2	9,737	9,914	99.93	99.95
3	9,747	9,295	99.96	99.64	4.6	3.4	3	9,911	9,959	99.98	99.98
4	9,883	9,721	99.98	99.80	6.0	4.4	4	9,937	9,963	99.99	99.98
5	9,932	9,874	99.99	99.90	7.6	5.5	5	9,948	9,966	99.99	99.98
6	9,936	9,913	99.99	99.94	9.1	6.6	6	9,948	9,966	99.98	99.98
7	9,935	9,936	99.99	99.96	10.6	7.7	7	9,941	9,974	99.98	99.97
8	9,939	9,948	99.99	99.97	12.2	8.9	8	9,948	9,975	99.98	99.97
9	9,944	9,955	99.98	99.97	13.7	10.0	9	9,954	9,979	99.98	99.98
10	9,956	9,958	99.99	99.98	15.2	11.2	10	9,960	9,983	99.98	99.98
15	9,972	9,975	99.99	99.98	23.0	16.8	15	9,973	9,985	99.98	99.98
20	9,955	9,958	99.98	99.98	30.7	22.1	20	9,958	9,976	99.98	99.98

Longer reads improve sensitivity, with assembled 42nt and 50nt reads identifying 82.0% and 94.7% of rare 1 p.p.m. clonotypes

## Acknowledgements

### References:

- Arestia PT et al. (1999) A direct estimate of the human T cell receptor diversity. *Science*. 286, 958-961.
- Gorski J et al. (1994) Circulating T cell repertoire complexity in normal individuals and bone marrow recipients analyzed by CDR3 size spectratyping. Correlation with immune status. *J. Immunol.* 152, 5109-5119.
- Pannetier C et al. (1993) The sizes of the CDR3 hypervariable regions of the murine T cell receptor beta chains vary as a function of the recombined germ-line segments. *Proc. Natl. Acad. Sci. USA*. 90, 4319-4323.
- Warren RL et al. (2009) Profiling model T cell metagenomes with short reads. *Bioinformatics*. doi:10.1093/bioinformatics/btp010
- Warren RL et al. (2009) Assembling millions of short DNA sequences with ISSAKE. *Bioinformatics*. 23:500-1 (epub dec09)

Funding: GenomeBritishColumbia