# Machine learning for the analysis of adaptive immune receptors and repertoires
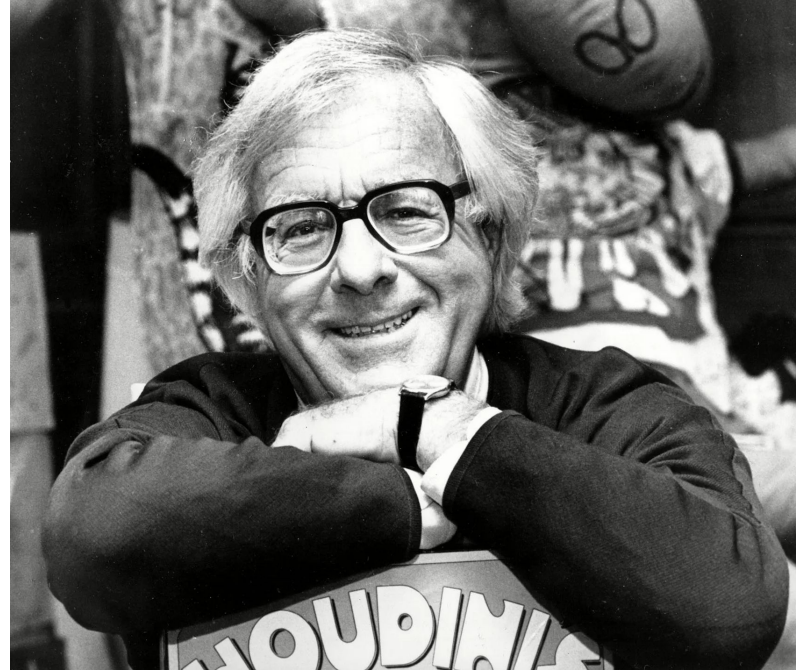
Maria Chernigovskaya
mariiac@uio.no

Milena Pavlović
milenpa@uio.no

AIRR Community Webinar

November 15, 2022

*"Life is trying things to see if they work."*
- Ray Bradbury

*Machine learning for the analysis of*
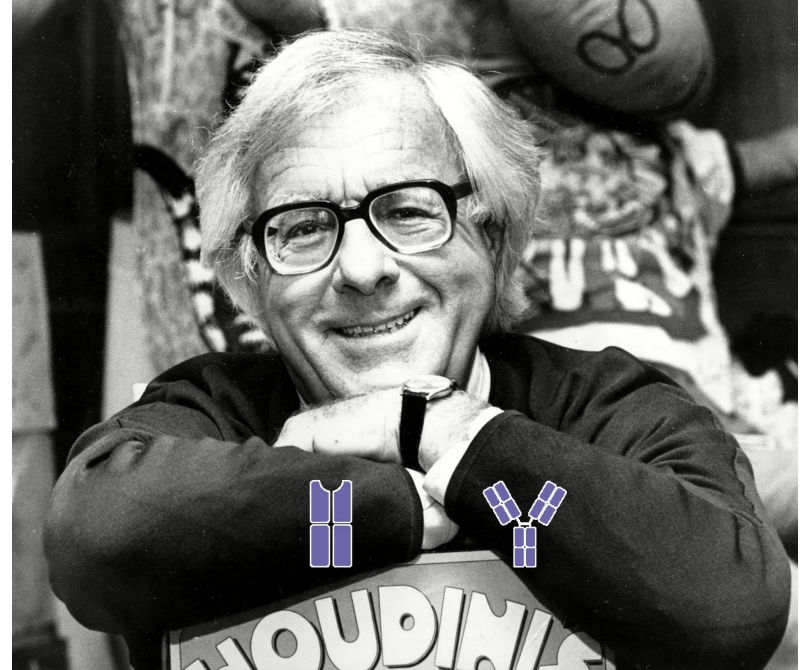*adaptive immune receptors and repertoires*

*"~~Life~~ is trying things to see if they work."*
                                    - AIRR researches

*Machine learning for the analysis of adaptive immune receptors and repertoires*

*"~~Life~~ is trying things to see if they work."*
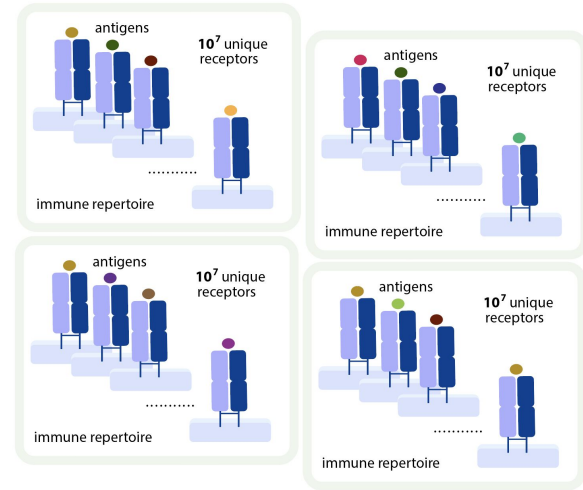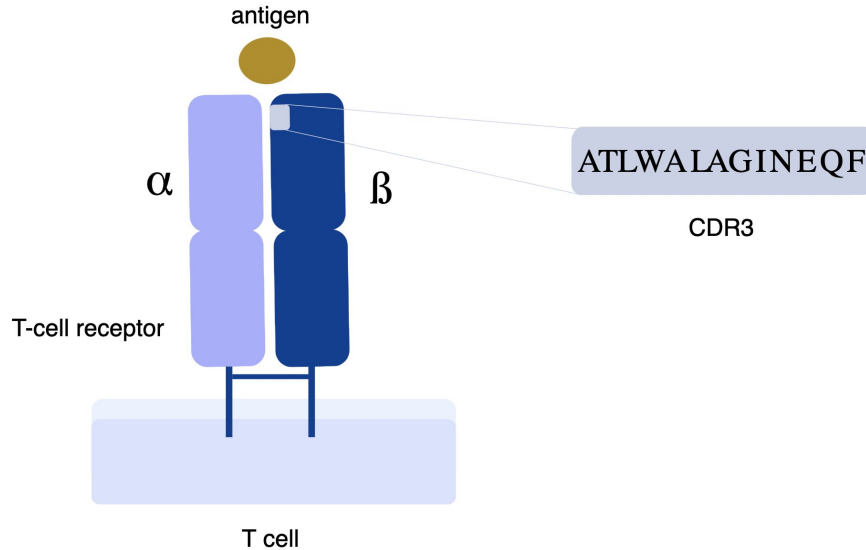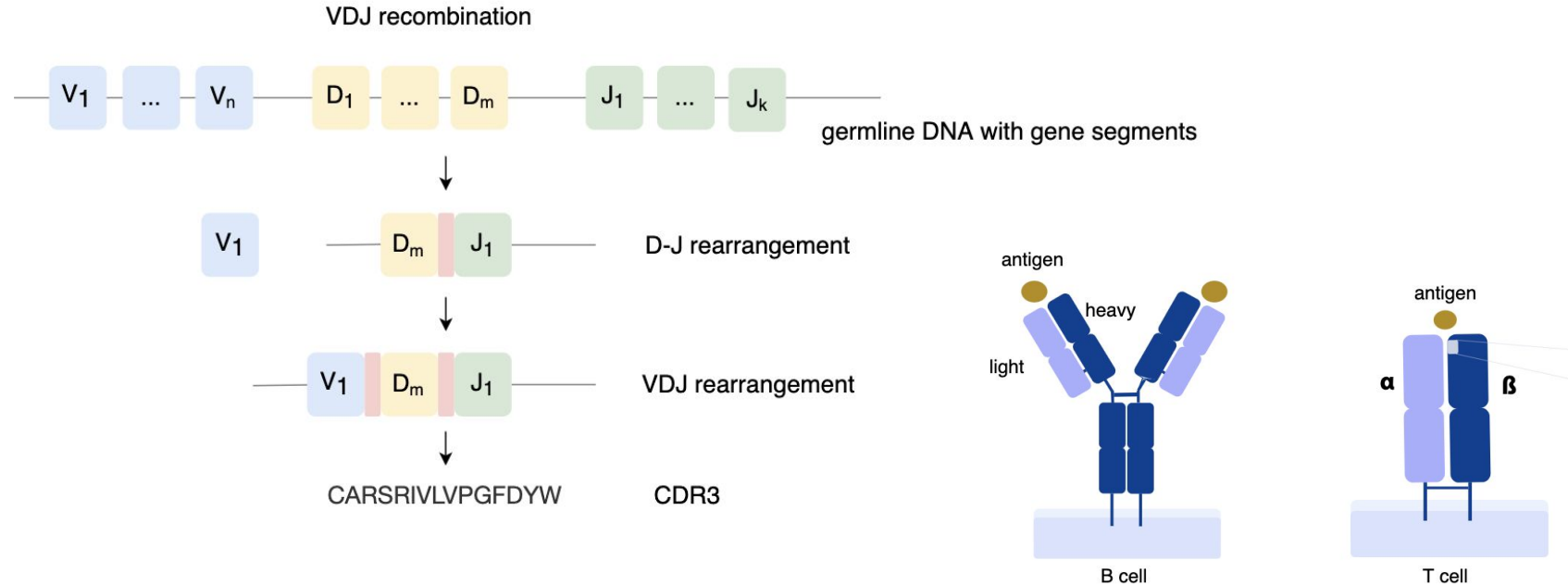*- AIRR researches*

Or not?

# Adaptive immune receptors (AIRs) and repertoires (AIRRs)

❏ Adaptive immune receptors (AIRs)

❏ Adaptive immune repertoires (AIRRs)

antigen

α ß

ATLWALAGINEQF

CDR3

T-cell receptor

T cell

antigens 10⁷ unique receptors

immune repertoire

antigens 10⁷ unique receptors

immune repertoire

antigens 10⁷ unique receptors

immune repertoire

antigens 10⁷ unique receptors

immune repertoire

# V(D)J recombination assembles AIRs (BCRs or TCRs)



VDJ recombination

V₁ ... Vₙ    D₁ ... Dₘ    J₁ ... Jₖ    germline DNA with gene segments

V₁    Dₘ J₁    D-J rearrangement

V₁ Dₘ J₁    VDJ rearrangement

CARSRIVLVPGFDYW    CDR3

antigen heavy light — B cell

antigen α β — T cell

Alt et al. 1980, 1984, 1992, Bassing et al. 2000, 2002, Bareto et al. 2000, Schatz et al. 2012

# Overview of AIRR data on the receptor level

TCR/BCR sequence   ⟶   TCR/BCR structure   ⟶   TCR/BCR function (binding)
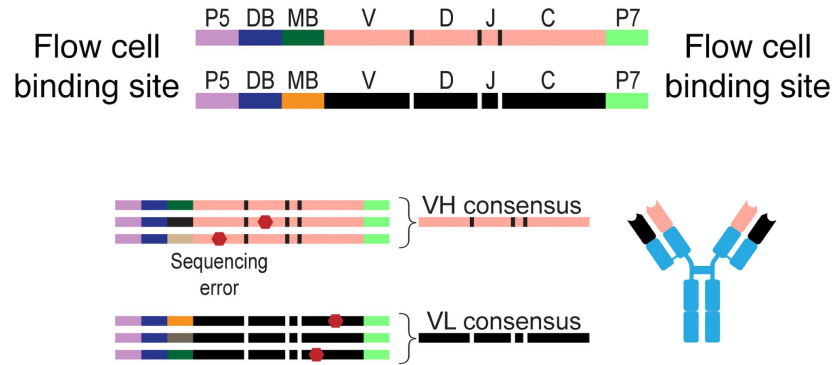
- ❏    Nt or aa
- ❏    Full-length or CDR3



5′UTR   L1   L2   FR1    V    D   J    FR3 CDR3 FR4

```
TRBV7-3 + CASSDRHQPQHF + TRBJ2-7
```



TCR constant domain

TCR variable domain

Peptide

MHC

β2-microglobulin

Bolotin et al. 2015, Zoete et al. 2013

# Overview of AIRR data on the repertoire level
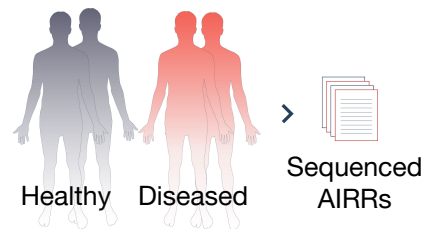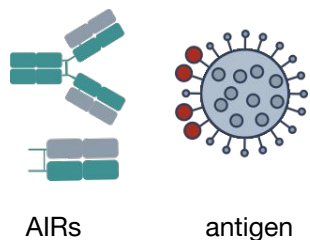
❏ AIR-seq (Bulk or Single cell, with or without UMIs)


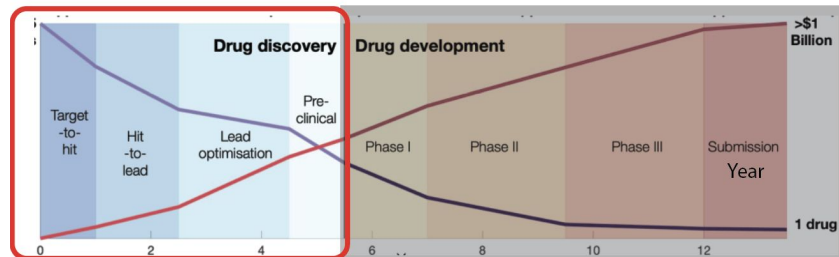
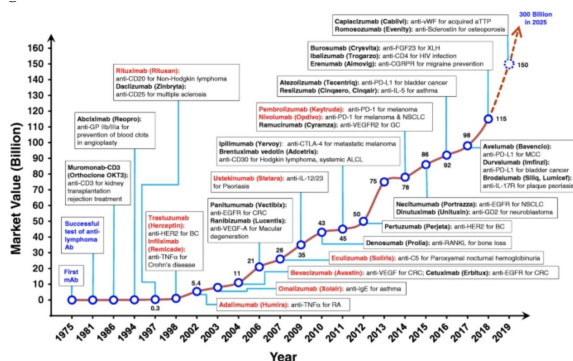Benichou et al. 2012, Yaari et al. 2015, Brown et al. 2019

❏ Antibody repertoire proteomics (Cheung et al. 2012, Sato et al. 2012, Wine et al. 2015, Snapkov et al. 2021)
❏ Paired AIRs mapped to Ag specificity (Setliff et al. 2019, "A new way of exploring immunity" 2020)
❏ Paired AIRs + gene expression (Tu et al. 2019, Mathew et al. 2021, Shlesinger et al. 2022, Gao et al. 2022, Stephenson et al. 2021)

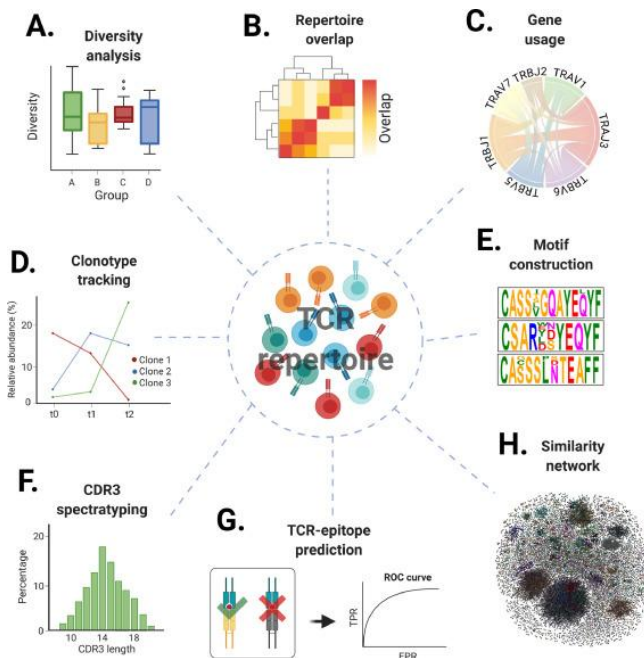# What immunological questions can we ask using AIRR data?

❏ Can we use an immune repertoire for disease diagnostics?



AIRs          antigen

Healthy   Diseased    Sequenced
                      AIRRs

❏ How can we improve vaccines/CAR-T/mAbs/other therapeutics design?



Lu et al. 2020

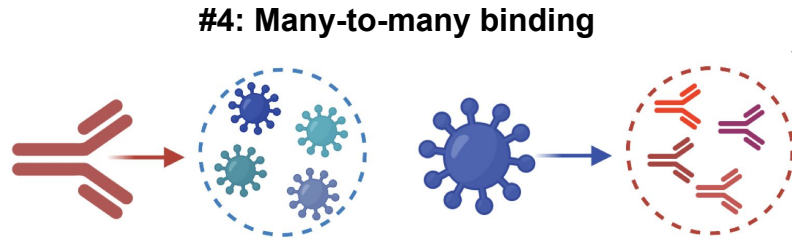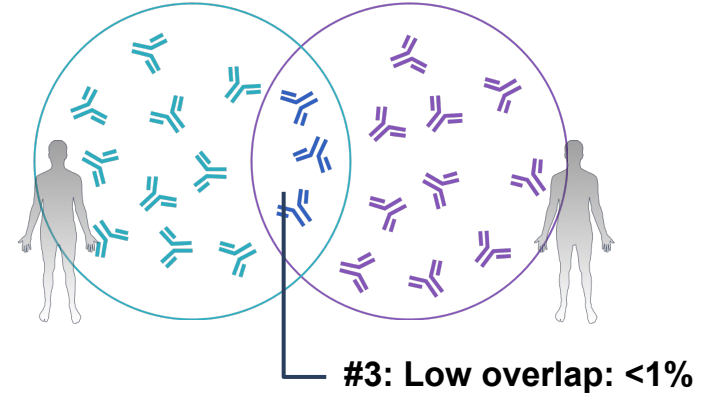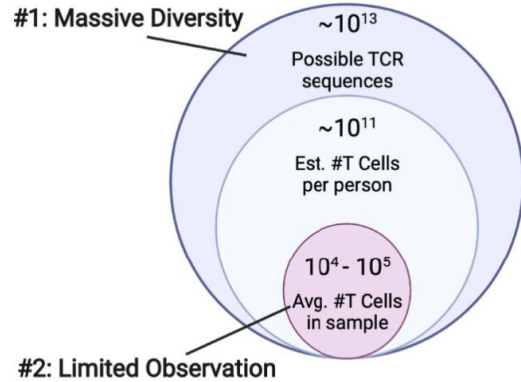# Some of these questions can be answered with various AIRR computational tools



Valkiers et al. 2022



Miho et al. 2018

# Challenges in computational analysis on AIRR data



Katayama et al. 2022, Greiff et al. 2020

Machine learning (ML) provides various approaches to detect signals in complex high-dimensional data

# What is machine learning?

❏ Machine learning (ML): a set of pattern recognition and function approximation techniques that find patterns within groups in (large amounts of) data
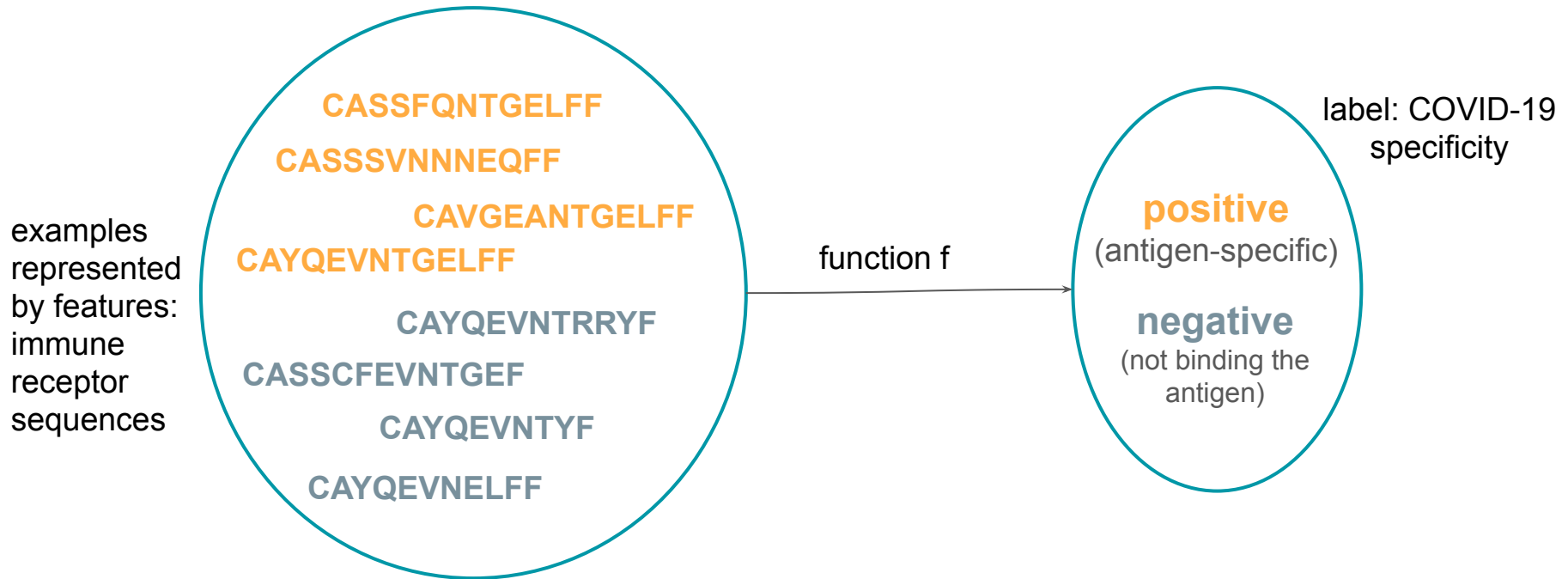
# What is machine learning?

❏ Machine learning (ML): a set of pattern recognition and function approximation techniques that find patterns within groups in (large amounts of) data
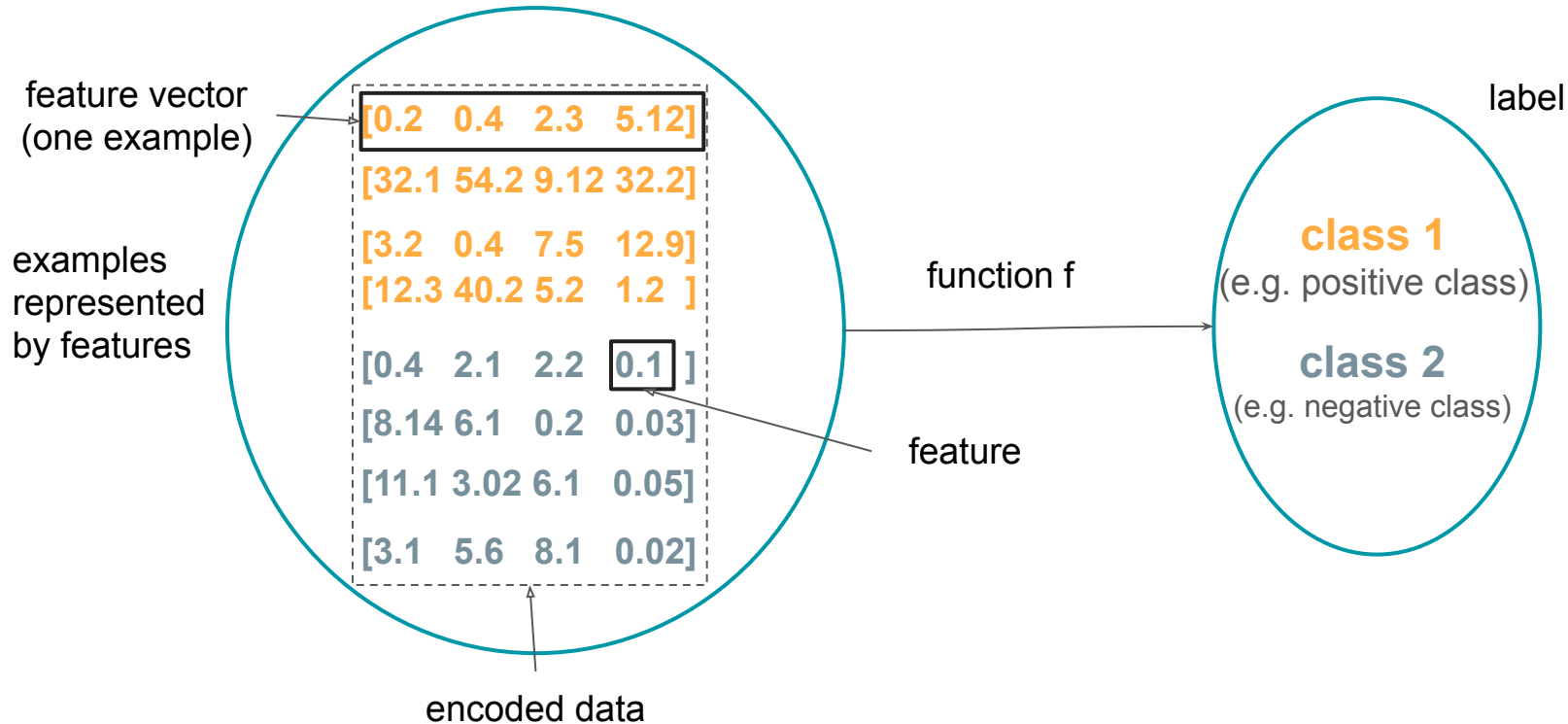
❏ A set of methods that allow for making inferences about the data

# Machine learning as a function approximation task

**CASSFQNTGELFF**

**CASSSVNNNEQFF**

**CAVGEANTGELFF**

**CAYQEVNTGELFF**

CAYQEVNTRRYF

CASSCFEVNTGEF

CAYQEVNTYF

CAYQEVNELFF

examples represented by features: immune receptor sequences

function f

label: COVID-19 specificity

**positive**
(antigen-specific)

**negative**
(not binding the antigen)

# Machine learning as a function approximation task

feature vector
(one example)

label

[0.2    0.4    2.3    5.12]

[32.1  54.2  9.12  32.2]

[3.2    0.4    7.5    12.9]
[12.3  40.2  5.2    1.2  ]

class 1
(e.g. positive class)

examples
represented
by features

function f

[0.4    2.1    2.2    0.1 ]

[8.14  6.1    0.2    0.03]

class 2
(e.g. negative class)

[11.1  3.02  6.1    0.05]

feature

[3.1    5.6    8.1    0.02]

encoded data

# Machine learning as a function approximation task

feature vector
(one example)

[0.2    0.4    2.3    5.12]

[32.1 54.2 9.12 32.2]

[3.2    0.4    7.5    12.9]
[12.3 40.2 5.2    1.2  ]

examples
represented
by features

[0.4    2.1    2.2    0.1 ]

[8.14 6.1    0.2    0.03]

[11.1 3.02 6.1    0.05]

[3.1    5.6    8.1    0.02]

Parameters of function
f are learned during
training

function f

feature

Representation
(encoding):
manually set or learned

encoded data

label

**class 1**
(e.g. positive class)

**class 2**
(e.g. negative class)

# Building predictive models

❏ Machine learning discovers statistical associations in the data
→ these associations enable good prediction

# Building predictive models

❏ Machine learning discovers statistical associations in the data
  → these associations enable good prediction

❏ Aim: get a good predictive model, but also get biological insight

  ❏ This is why we want the ML models to be interpretable

# Building predictive models

❏ Machine learning discovers statistical associations in the data
→ these associations enable good prediction

❏ Aim: get a good predictive model, but also get biological insight

  ❏ This is why we want the ML models to be interpretable

❏ Not causal relations, but starting points for further analyses

# There is a surge in (AIRR) ML studies

**Immunosequencing identifies signatures of cytomegalovirus exposure history and HLA-mediated effects on the T cell repertoire**

Ryan O Emerson, William S DeWitt, Marissa Vignali, Jenna Gravley, Joyce K Hu, Edward J Osborne, Cindy Desmarais, Mark Klinger, Christopher S Carlson, John A Hansen, Mark Rieder & Harlan S Robins

**Predicting unseen antibodies' neutralizability via adaptive graph neural networks**

Jie Zhang, Yishan Du, Pengfei Zhou, Jinru Ding, Shuai Xia, Qian Wang, Feiyang Chen, Mu Zhou, Xuemei Zhang, Weifeng Wang, Hongyan Wu, Lu Lu & Shaoting Zhang

**DeepTCR: a deep learning framework for understanding T-cell receptor sequence signatures within complex T-cell repertoires**

John-William Sidhom, H. Benjamin Larman, Petra Ross-MacDonald, Megan Wind-Rotolo, Drew M. Pardoll, Alexander S. Baras

**TITAN: T-cell receptor specificity prediction with bimodal attention networks**

Anna Weber, Jannis Born, María Rodriguez Martínez

**Biophysicochemical Motifs in T-cell Receptor Sequences Distinguish Repertoires from Tumor-Infiltrating Lymphocyte and Adjacent Healthy Tissue**

Jared Ostmeyer, Scott Christley, Inimary T. Toby, and Lindsay G. Cowell

**Predicting antigen specificity of single T cells based on TCR CDR3 regions**

David S Fischer, Yihan Wu, Benjamin Schubert, Fabian J Theis
Author Information

**Attentive Cross-Modal Paratope Prediction**

Andreea Deac, Petar Veličković, and Pietro Sormanni

**Parapred: antibody paratope prediction using convolutional and recurrent neural networks**

Edgar Liberis, Petar Veličković, Pietro Sormanni, Michele Vendruscolo, Pietro Liò

**Mining adaptive immune receptor repertoires for biological and clinical information using machine learning**

Victor Greiff[1], Gur Yaari[2], Lindsay G. Cowell[3]

**Modern Hopfield Networks and Attention for Immune Repertoire Classification**

Authors
Michael Widrich, Bernhard Schäfl, Milena Pavlović, Hubert Ramsauer, Lukas Gruber, Markus Holzleitner, Johannes Brandstetter, Geir Kjetil Sandve, Victor Greiff, Sepp Hochreiter, Günter Klambauer

**Capturing the differences between humoral immunity in the normal and tumor environments from repertoire-seq of B-cell receptors using supervised machine learning**

Hiroki Konishi, Daisuke Komura, Hiroto Katoh, Shinichiro Atsumi, Hirotomo Koda, Asami Yamamoto, Yasuyuki Seto, Masashi Fukayama, Rui Yamaguchi, Seiya Imoto & Shumpei Ishikawa

**Detection of Enriched T Cell Epitope Specificity in Full T Cell Receptor Sequence Repertoires**

Sofie Gielis[1,2,3], Pieter Moris[1,3†], Wout Bittremieux[1,3†], Nicolas De Neuter[1,2,3], Benson Ogunjimi[2,5,6,7], Kris Laukens[1,2,3‡] and Pieter Meysman[1,2,3*‡]

**De novo prediction of cancer-associated T cell receptors for noninvasive cancer detection**

Daria Beshnova[1], Jianfeng Ye[1], Oreoluwa Onabolu[2], Benjamin Moon[3], Wenxin Zheng[4], Yang-Xin Fu[3,5], James Brugarolas[2], Jayanthi Lea[4] and Bo Li[1,5,*]

[1]Lyda Hill Department of Bioinformatics, UT Southwestern Medical Center, Dallas, TX 75390, USA.
[2]Department of Internal Medicine, UT Southwestern Medical Center, Dallas, TX 75390, USA.
[3]Department of Pathology, UT Southwestern Medical Center, Dallas, TX 75390, USA.
[4]Department of Obstetrics and Gynecology, UT Southwestern Medical Center, Dallas, TX 75390, USA.
[5]Department of Immunology, UT Southwestern Medical Center, Dallas, TX 75390, USA.
↵*Corresponding author. Email: bo.li@utsouthwestern.edu
← Hide authors and affiliations

**+ Many more!**

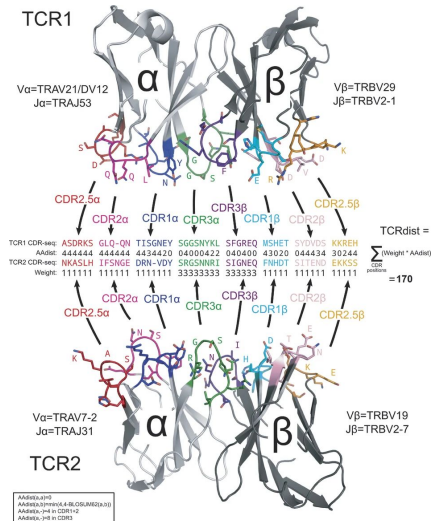# ML application areas in AIRR analyses

# Predicting receptor specificity

❏ Examining sequence similarity:

Find the similarity between known positive and negative sequences, and predict the specificity to be the same as the sequences in the closest proximity
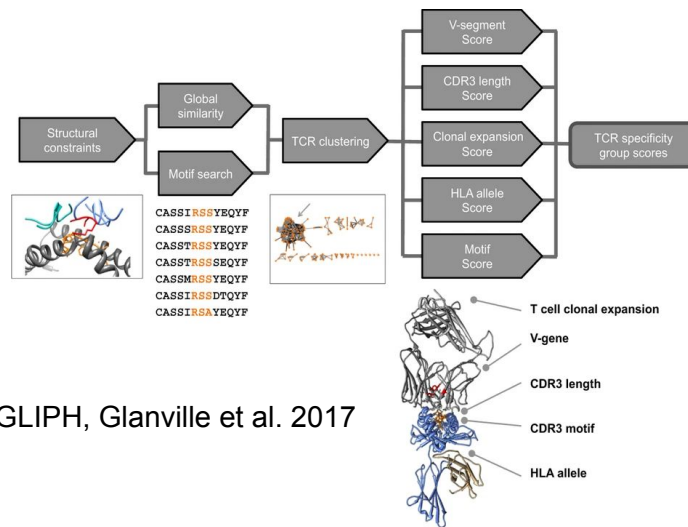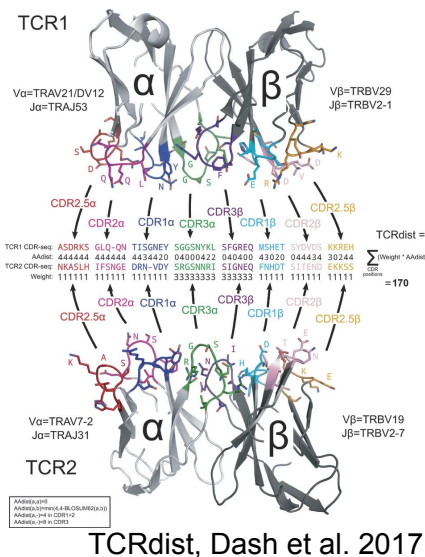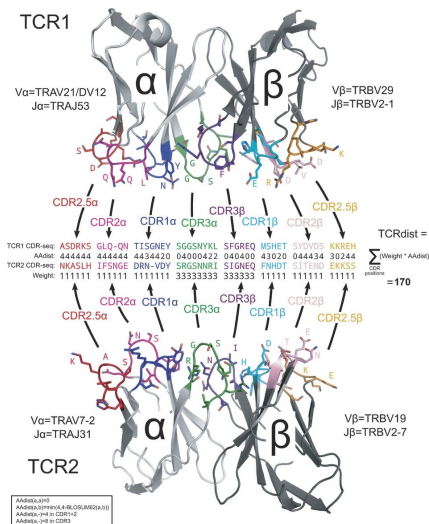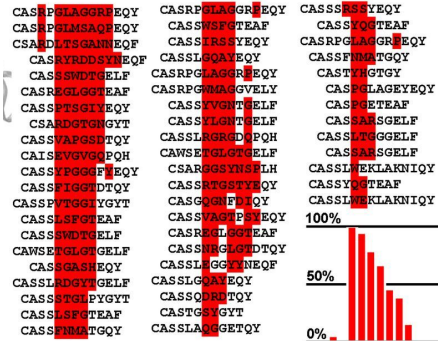
# Predicting receptor specificity

❏ Examining sequence similarity:

Find the similarity between known positive and negative sequences, and predict the specificity to be the same as the sequences in the closest proximity



TCRdist, Dash et al. 2017

# Predicting receptor specificity

❏ Examining sequence similarity:

Find the similarity between known positive and negative sequences, and predict the specificity to be the same as the sequences in the closest proximity



TCRdist, Dash et al. 2017

GLIPH, Glanville et al. 2017

# Predicting receptor specificity

❏ Examining sequence similarity:

Find the similarity between known positive and negative sequences, and predict the specificity to be the same as the sequences in the closest proximity



TCRdist, Dash et al. 2017

GLIPH, Glanville et al. 2017

iSMART, Zhang et al. 2020

# Predicting receptor specificity

❏ Discovering short motifs in the sequence that are indicative of its specificity

❏ Predictions made based on physicochemical properties of receptors



Ostmeyer et al. 2019

# Predicting receptor specificity

❏ Discovering short motifs in the sequence that are indicative of its specificity

❏ Predictions made based on physicochemical properties of receptors



Ostmeyer et al. 2019

Chronister et al. 2021

# Predicting receptor specificity

❏ Discovering short motifs in the sequence that are indicative of its specificity

❏ Predictions made based on physicochemical properties of receptors

Ostmeyer et al. 2019

Chronister et al. 2021

Kanduri et al. 2022

# Predicting receptor specificity

❏   Modeling antibody-antigen interactions



Zhang et al. 2022

# Predicting receptor specificity

❏ Modeling antibody-antigen interactions

❏ Using structural information



Zhang et al. 2022

**DeepAIR: a deep-learning framework for effective integration of sequence and 3D structure to enable adaptive immune receptor analysis**

Yu Zhao, 🆔 Bing He, 🆔 Chen Li, Zhimeng Xu, Xiaona Su, Jamie Rossjohn, 🆔 Jiangning Song, Jianhua Yao

**doi:** https://doi.org/10.1101/2022.09.30.510251

**DLAB: deep learning methods for structure–based virtual screening of antibodies** 🔓

Constantin Schneider, Andrew Buchanan, Bruck Taddese, Charlotte M Deane ✉

*Bioinformatics*, Volume 38, Issue 2, 15 January 2022, Pages 377–383, https://doi.org/10.1093/bioinformatics/btab660

**Published:** 21 September 2021   **Article history** ▾

31

# Analysis of AIR data

❑ Learning a latent representation using sequence and gene expression data

# Analysis of AIR data

❏ Learning a latent representation using sequence and gene expression data



mvTCR, Drost et al. 2022

# Analysis of AIR data

❏ Learning a latent representation using sequence and gene expression data



mvTCR, Drost et al. 2022



Benisse, Zhang et al. 2022

# Analysis of AIR data

❏ Learning a latent representation using sequence and gene expression data
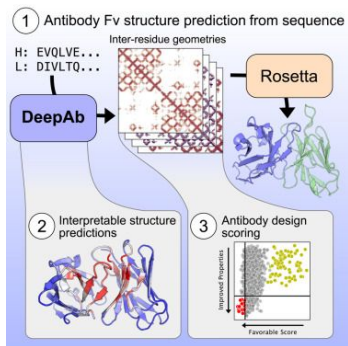


Benisse, Zhang et al. 2022



mvTCR, Drost et al. 2022



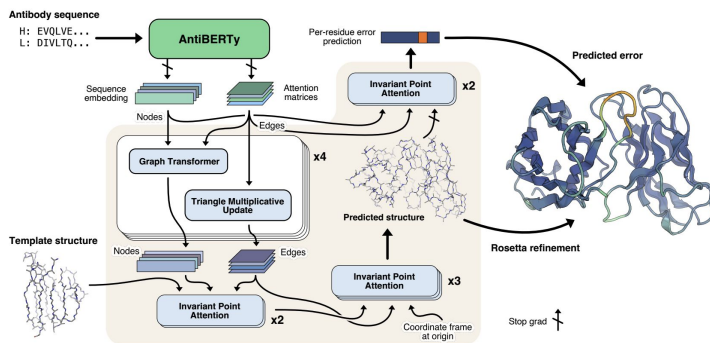GIANA, Zhang et al. 2021

# Predicting the 3D structure of AIRs

❏ Antibody-specific methods achieve better prediction performance than generic protein structure prediction tools
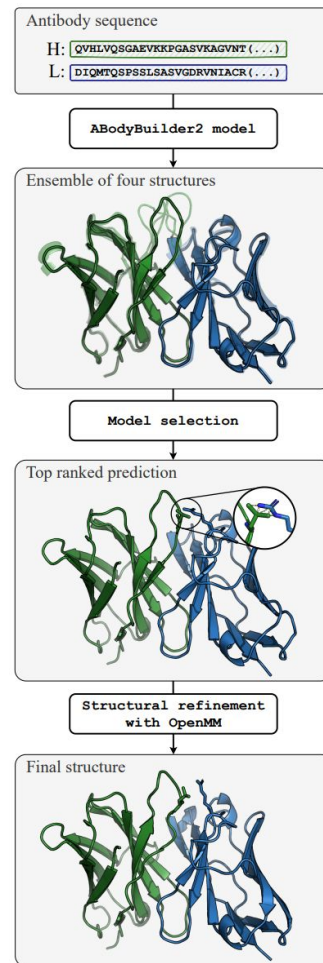


ABlooper, Abanades et al. 2022

# Predicting the 3D structure of AIRs

❏ Antibody-specific methods achieve better prediction
   performance than generic protein structure prediction tools



ABlooper, Abanades et al. 2022



DeepAb, Ruffolo et al. 2022

# Predicting the 3D structure of AIRs

❏ Antibody-specific methods achieve better prediction performance than generic protein structure prediction tools



ABlooper, Abanades et al. 2022



DeepAb, Ruffolo et al. 2022

IgFold, Ruffolo et al. 2022

# Predicting the 3D structure of AIRs

❏ Antibody-specific methods achieve better prediction performance than generic protein structure prediction tools



ABlooper, Abanades et al. 2022



DeepAb, Ruffolo et al. 2022



IgFold, Ruffolo et al. 2022



ImmuneBuilder, Abanades et al. 2022

# Language models for antibody sequences

❏   Some of the previous models are based
    on language models

# Language models for antibody sequences

❏ Some of the previous models are based on language models

❏ Necessary to formalize the "antibody language"

# Language models for antibody sequences

❏ Some of the previous models are based on language models

❏ Necessary to formalize the "antibody language"



ImmunoLingo, Vu et al. 2022

# Language models for antibody sequences

- ❏ Some of the previous models are based on language models

- ❏ Necessary to formalize the "antibody language"

- ❏ Improved interpretability through formalization



ImmunoLingo, Vu et al. 2022

# Language models for antibody sequences

❏ Some of the previous models are based on language models

❏ Necessary to formalize the "antibody language"

❏ Improved interpretability through formalization

❏ Potential aim: therapeutics design



ImmunoLingo, Vu et al. 2022

# Generative models for AIRs

# Generative models for AIRs

❏ Modeling the VDJ recombination process (naive AIRs)
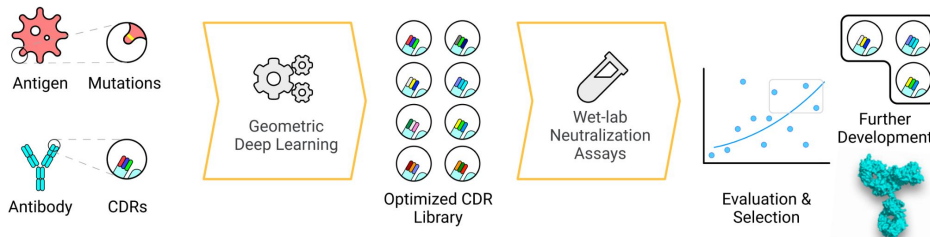
IGoR, Marcou et al. 2018,
doi:10.1038/s41467-018-02832-w

Davidsen et al. 2019, doi: 10.7554/eLife.46935

# Generative models for AIRs

❏ Modeling the VDJ recombination process (naive AIRs)



IGoR, Marcou et al. 2018,
doi:10.1038/s41467-018-02832-w

Davidsen et al. 2019, doi: 10.7554/eLife.46935

❏ Modeling antigen–specific antibodies directly



Saka et al. 2021,
doi: 10.1038/s41598-021-85274-7

Shan et al. 2022,
doi: 10.1073/pnas.2122954119

# Antibody design with machine learning

- ❏ Epitope specificity, affinity and developability

- ❏ Public repositories: iReceptor, IEDB, AbDb, AgAbDb

- ❏ Synthetic data: Absolut!

48

# TCRs and peptide-MHC complexes

❏    For TCRs to recognize a peptide, it has to be presented by the MHC complex

# TCRs and peptide-MHC complexes

❏ For TCRs to recognize a peptide, it has to be presented by the MHC complex

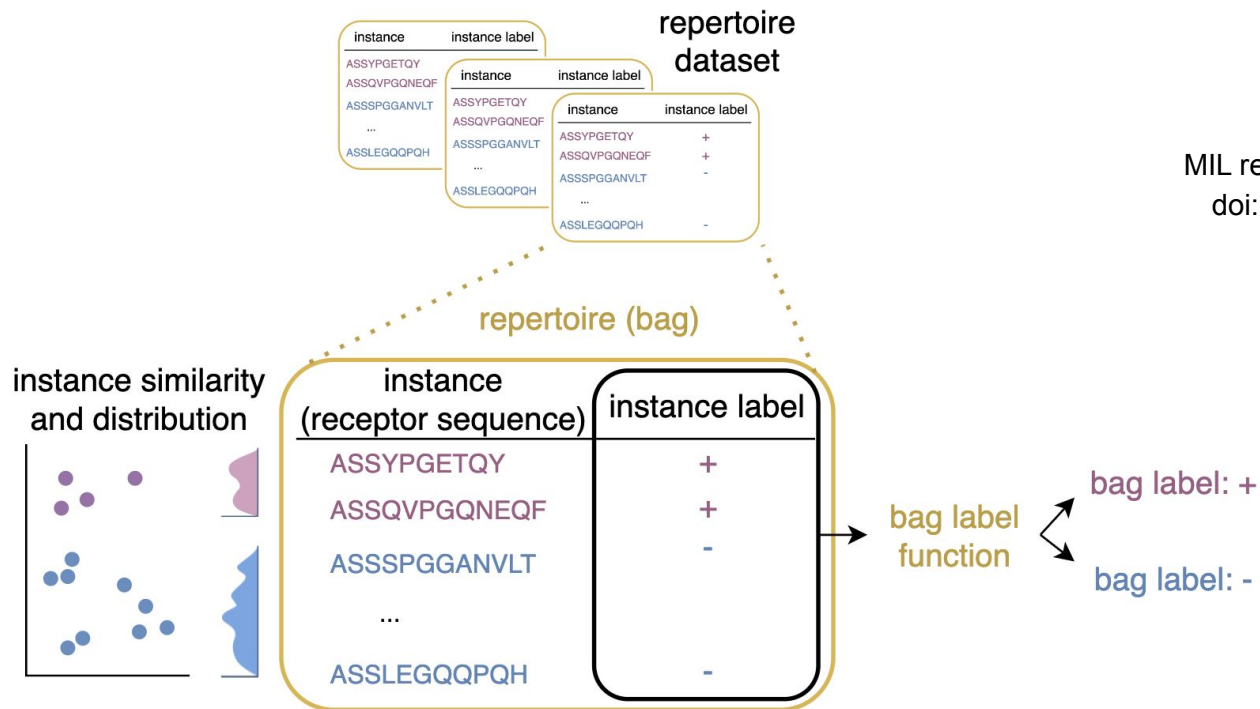❏ Tasks: predicting peptide-MHC binding, predicting the binding of a TCR to pMHC complex

# TCRs and peptide-MHC complexes

❏ For TCRs to recognize a peptide, it has to be presented by the MHC complex

❏ Tasks: predicting peptide-MHC binding, predicting the binding of a TCR to pMHC complex



NNAlign in the review by Nielsen et al. 2020,
doi:10.1146/annurev-biodatasci-021920-100259

51

# Diagnosing immune-related diseases with AIRRs

❏ Repertoire classification is a multiple instance learning (MIL) problem

# Diagnosing immune-related diseases with AIRRs

❏ Repertoire classification is a multiple instance learning (MIL) problem



MIL review: Carbonneau et al. 2018
doi:10.1016/j.patcog.2017.10.009

# Custom ML approaches for diagnostics

## Immunosequencing identifies signatures of cytomegalovirus exposure history and HLA-mediated effects on the T cell repertoire

Ryan O Emerson ✉, William S DeWitt, Marissa Vignali, Jenna Gravley, Joyce K Hu, Edward J Osborne, Cindy Desmarais, Mark Klinger, Christopher S Carlson, John A Hansen, Mark Rieder & Harlan S Robins

---

## DeepTCR is a deep learning framework for revealing sequence concepts within T-cell repertoires

John-William Sidhom ✉, H. Benjamin Larman, Drew M. Pardoll & Alexander S. Baras

---

## Multiple Instance Neural Networks Based on Sparse Attention for Cancer Detection using T−cell Receptor Sequences

Younghoon Kim, Tao Wang, Danyi Xiong, Xinlei Wang, Seongoh Park

## Modern Hopfield Networks and Attention for Immune Repertoire Classification

Michael Widrich, Bernhard Schäfl, Hubert Ramsauer, Milena Pavlović, Lukas Gruber, Markus Holzleitner, Johannes Brandstetter, Geir Kjetil Sandve, Victor Greiff, Sepp Hochreiter, Günter Klambauer

## Biophysicochemical Motifs in T-cell Receptor Sequences Distinguish Repertoires from Tumor-Infiltrating Lymphocyte and Adjacent Healthy Tissue FREE

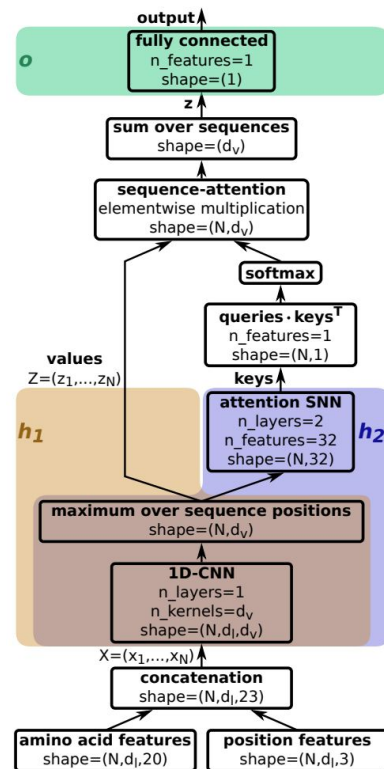Jared Ostmeyer ⓘ; Scott Christley; Inimary T. Toby; Lindsay G. Cowell ✉

## Disease diagnostics using machine learning of immune receptors

ⓘ Maxim E. Zaslavsky, ⓘ Nikhil Ram-Mohan, ⓘ Joel M. Guthridge, ⓘ Joan T. Merrill, ⓘ Jason D. Goldman, ⓘ Ji-Yeun Lee, ⓘ Krishna M. Roskin, ⓘ Charlotte Cunningham-Rundles, ⓘ M. Anthony Moody, ⓘ Barton F. Haynes, ⓘ Benjamin A. Pinsky, ⓘ James R. Heath, ⓘ Judith A. James, ⓘ Samuel Yang, ⓘ Catherine A. Blish, ⓘ Robert Tibshirani, ⓘ Anshul Kundaje, ⓘ Scott D. Boyd
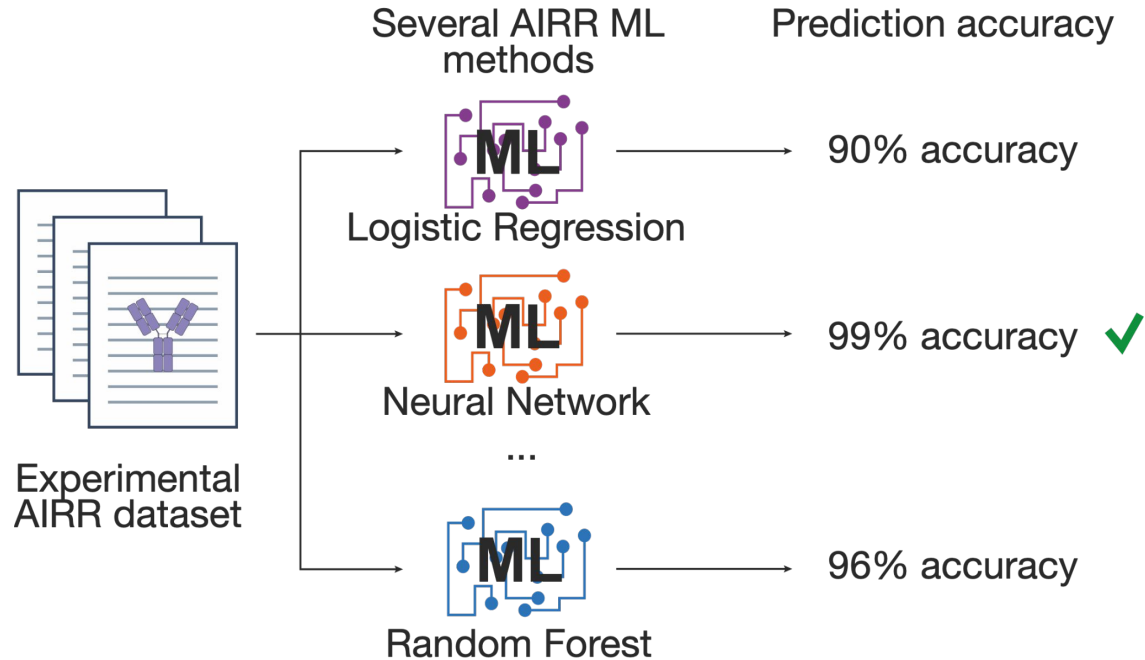
DeepRC, Widrich et al. 2020

Different ML methods have different underlying assumptions
– those should be conscious choices to reflect the problem domain

How do we ensure that the method can be applied to unseen receptors or repertoires?
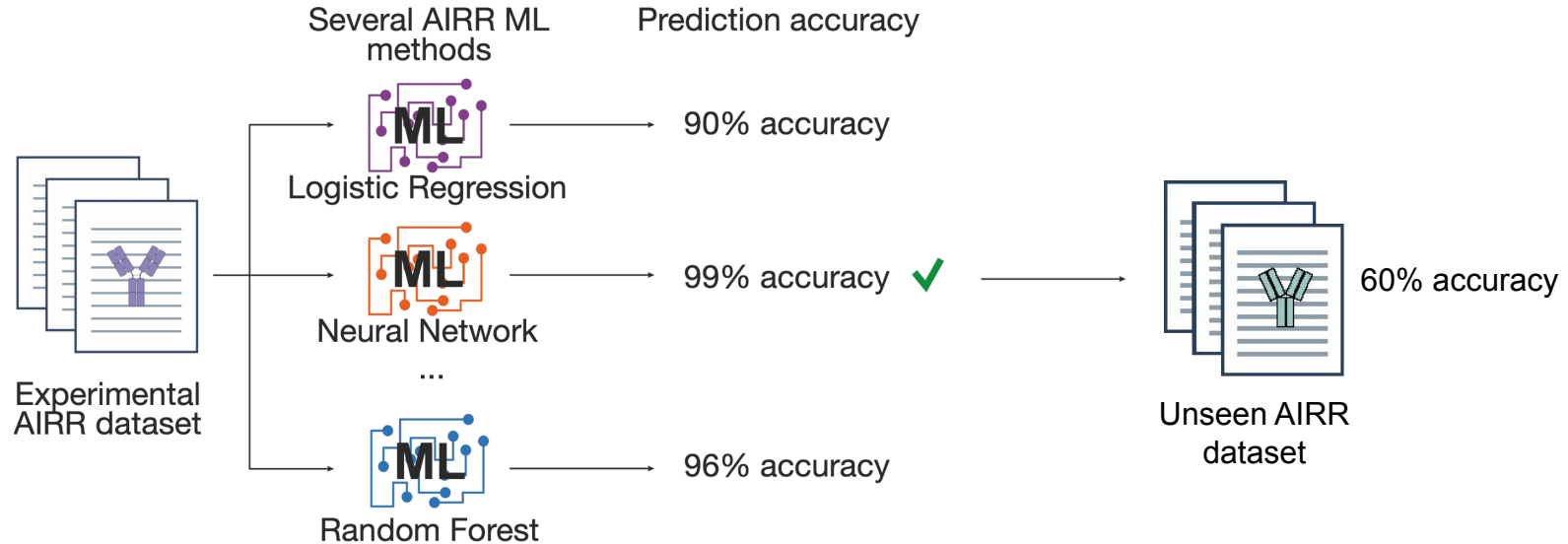
(generalizability of ML methods)

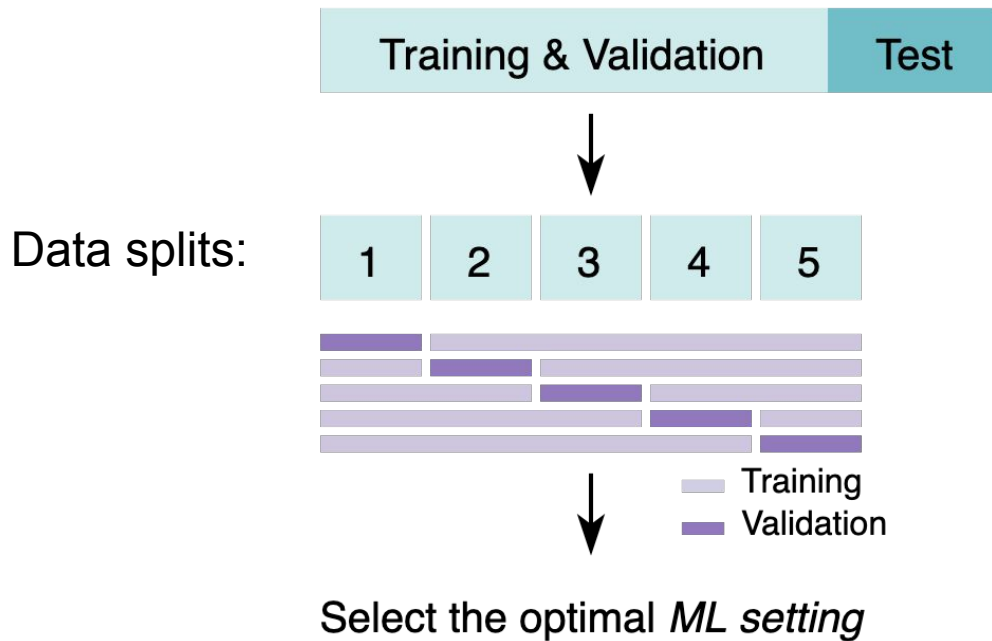# A naive way to perform AIRR ML

# Experimental data challenges in method development

- ❏ Usually small dataset size
    - ❏ One of the largest AIRR studies: Liu et al., 2019 — 877 repertoires,18 million unique TCRs, Snyder et al., 2020 — 1815 covid TCR repertoires + 3500 controls
    - ❏ Is the dataset representative? Performance estimation problems

- ❏ Available only for one particular problem setup
    - ❏ What if the data was a bit different? Sensitivity estimation problem

- ❏ No ground truth information
    - ❏ What was learned? Generalizability problem

# Will our ML method also work good on an unseen data?

# Nested cross validation might improve generalizability

# Nested cross validation might improve generalizability

# Nested cross validation might improve generalizability

# There are several levels of ML verification



Can we generalise?

**Model Verification**

Train Data → ML Model → Prediction

Test Data → ML Model → Prediction ✔

**Explanation Verification**

ML Model → Explanation → Hypothesis

Synthetic Data → Synthetic ground truth

**Knowledge Verification**

ML Model → Explanation → Hypothesis

Data → Lab Experiment → Well-studied mechanism

If well-studied prior knowledge is **not** available

If well-studied prior knowledge is available

# What can be a prior knowledge in AIRR case?



Ground truth (unknown for experimental data)

Model construction → AIRR generation model → Realisation → AIRR dataset

$\Theta$ AIRR model parameters

VDJ recombination

$V_1$ ... $V_n$  $D_1$ ... $D_m$  $J_1$ ... $J_k$  germline DNA with gene segments

$V_1$  $D_m$ $J_1$  D-J rearrangement

$V_1$ $D_m$ $J_1$  VDJ rearrangement

Several AIRR ML methods

Experimental AIRR dataset

Prediction accuracy

**ML** Logistic Regression → 90% accuracy

**ML** Neural Network → 99% accuracy ✔

...

**ML** Random Forest → 96% accuracy

# AIRR ML methods should be also benchmarked on ground truth synthetic AIRR data



*Can we identify the parameters?*

- ❏ High accuracy
- ❏ Have we learned the ground truth?

# Current VDJ simulation frameworks have pros and cons

❏ IGoR (nt) / OLGA (aa) (Marcou et al. 2017)



$$P(\text{scenario}) = P(V)P(J|V)\,P(D|V,J)P(\text{del}V|V)$$
$$\times\ P(\text{del}J|J)\,P(\text{del}D5'|D)\,P(\text{del}D3'|\text{del}D5',D)$$
$$\times\ P(\text{insVD})\prod_{i}^{\text{InsVD}} P(n_i|n_{i-1})$$
$$\times\ P(\text{insDJ})\prod_{i}^{\text{InsDJ}} P(n_i|n_{i-1})$$

+   Accurate VDJ recombination model
+   Generation probability evaluation
+   Fast (generates 100k seqs in 5 min)
-   No signal embedding
-   One dataset per one run

❏ immuneSIM (Weber et al. 2020)



+   Basic signal implantation (gapped k-mers)
+   Simulates clonal abundances
+   Productive receptors
-   Slow (100k seqs in 1 hour)
-   No generation model and generation probabilities
-   One dataset per one run

66

# Profiling AIRR ML models on a range of basic datasets

- ❑ OLGA TCRβ CDR3 sequences (aa) + implanted gapped 2—5-mers

- ❑ Identified parameter boundaries where baseline methods (Logistic Regression) already achieve high accuracy

Immune signal can be more complex!



Kanduri et al. 2022, doi: 10.1093/gigascience/giac046

67

# Definition of immune event and immune signal

We hypothesise that immune signal should be a substring of the receptor:

❏  (Gapped) k-mer (Akbar et al., 2021)
❏  Full-length receptor (Emerson et al., 2017)
❏  Motif (PWM with a fixed length)
❏  The most general definition: immune signal is a function: AIR → True/False



Chernigovskaya, unpublished

# A universal AIRR simulator wishlist

❏ Challenges in AIRR simulator development



*Synthetic AIRR data nativeness*

*Overlapping signals*

*Introducing simulation artifacts*

❏ Properties of a universal AIRR simulator



*Similar biological statistics distributions*

*Does not break biological properties of AIRR data*

*Can simulate signal of varying complexity*

Chernigovskaya, unpublished

# Framework for simulating a "native"-like AIR(R) datasets



Chernigovskaya, unpublished

Robert et al. 2022

Ground truth (unknown for experimental data)

Can we identify the parameters?

99% accuracy

- ❏ High accuracy
- ❏ Have we learned the ground truth?

# Conceptual problem: reproducible AIRR ML



How can we make all these studies reproducible?

# Recommendations for ML in biology

## DOME: recommendations for supervised machine learning validation in biology

Ian Walsh, Dmytro Fishman, Dario Garcia-Gasulla, Tiina Titma, Gianluca Pollastri, ELIXIR Machine Learning Focus Group, Jennifer Harrow ✉, Fotis E. Psomopoulos ✉ & Silvio C. E. Tosatto ✉

**Table 1 | Supervised ML in biology: concerns, the consequences they impart and recommendations**

| Broad topic | Be on the lookout for | Consequences | Recommendation(s) |
|---|---|---|---|
| Data | • Inadequate data size & quality<br>• Inappropriate partitioning, dependence between train and test data<br>• Class imbalance<br>• No access to data | • Data not representative of domain application<br>• Unreliable or biased performance evaluation<br>• Cannot check data credibility | • **Use independent optimization (training) and evaluation (testing) sets**. This is especially important for meta algorithms, where independence of multiple training sets must be shown to be independent of the evaluation (testing) sets.<br>• **Release data, preferably using appropriate long-term repositories, and include exact splits**.<br>• Offer sufficient evidence of data size & distribution being representative of the domain. |
| Optimization | • Overfitting, underfitting and illegal parameter tuning<br>• Imprecise parameters and protocols given | • Reported performance is too optimistic or too pessimistic<br>• The model models noise or misses relevant relationships<br>• Results are not reproducible | • **Clarify that evaluation sets were not used for feature selection, preprocessing steps or parameter tuning**.<br>• **Report indicators on training and testing data that can aid in assessing the possibility of under- or overfitting; for example, train vs. test error**.<br>• **Release definitions of all algorithmic hyperparameters, regularization protocols, parameters and optimization protocol**.<br>• For neural networks, release definitions of training and learning curves.<br>• Include explicit model validation techniques, such as $N$-fold cross-validation. |
| Model | • Unclear if black box or interpretable model<br>• No access to resulting source code, trained models & data<br>• Execution time impractical | • An interpretable model shows no explainable behavior<br>• Cannot cross compare methods & reproducibility, or check data credibility<br>• Model takes too much time to produce results | • **Describe the choice of black box or interpretable model. If interpretable, show examples of interpretable output**.<br>• Release documented source code + models + executable + user interface/webserver + software containers.<br>• Report execution time averaged across many repeats. If computationally tough, compare to similar methods. |
| Evaluation | • Performance measures inadequate<br>• No comparisons to baselines or other methods<br>• Highly variable performance | • Biased performance measures reported<br>• The method is falsely claimed as state-of-the-art<br>• Unpredictable performance in production | • **Compare with public methods & simple models (baselines)**.<br>• **Adopt community-validated measures and benchmark datasets for evaluation**.<br>• Compare related methods and alternatives on the same dataset.<br>• Evaluate performance on a final independent held-out set.<br>• **Use confidence intervals/error intervals and statistical tests to gauge prediction robustness**. |

Key recommendations are bolded.

# immuneML is a platform for development and transparent comparative evaluation of AIRR-ML methods
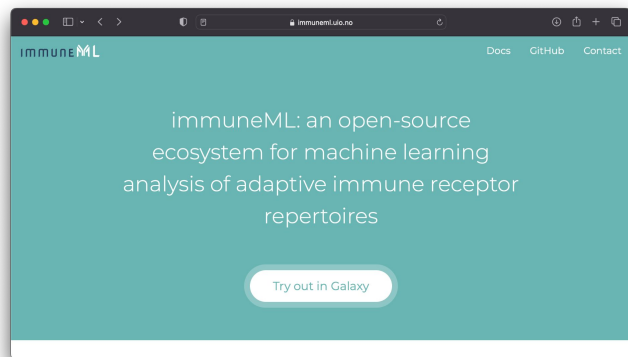


https://immuneml.uio.no

# Diagnosing diseases with AIRRs

distribution of
variables of interest

classes: diseased
and healthy

source population

# Diagnosing diseases with AIRRs



distribution of variables of interest

classes: diseased and healthy

source population

selecting participants for the study

study cohort

AIRR

$10^8$ receptors

# Diagnosing diseases with AIRRs

# Diagnosing diseases with AIRRs
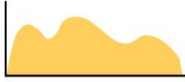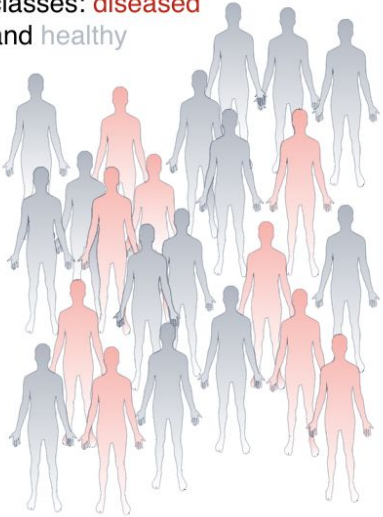
# Diagnosing diseases with AIRRs

# Diagnosing diseases with AIRRs



distribution of variables of interest → dataset shifts (changes in the distribution of variables of interest) → distribution of variables of interest

classes: diseased and healthy

AIRR $10^8$ receptors

~$10^6$ observed receptors per AIRR

batch 1  batch 2
batch 3  batch n

ASSISAG
ASSLEGQQ
ASSPGD

AIRR encoding and ML settings

source population — selecting participants for the study — sample collection and sequencing — sequenced AIRRs — fitting an ML model — application

# The causal inference framework

❏ Formally describes the data-generating process to discover causal effects between the variables in the process, under a set of assumptions (Pearl 2009)

❏ Causal effect of X on Y: the difference in the value of Y while changing X and keeping all other variables and conditions the same

Causality doesn't matter [too much] for prediction tasks, but when obtaining the data or applying methods to new populations, causality can help formalize and solve challenges even in predictive settings

# The causal inference framework in the AIRR field

❏ A causal model for a viral infection
(different for different diseases)



Pavlovic et al. 2022

# The causal inference framework in the AIRR field

- ❏ A causal model for a viral infection (different for different diseases)

- ❏ Selection bias:

  - ❏ preferential selection of study participants
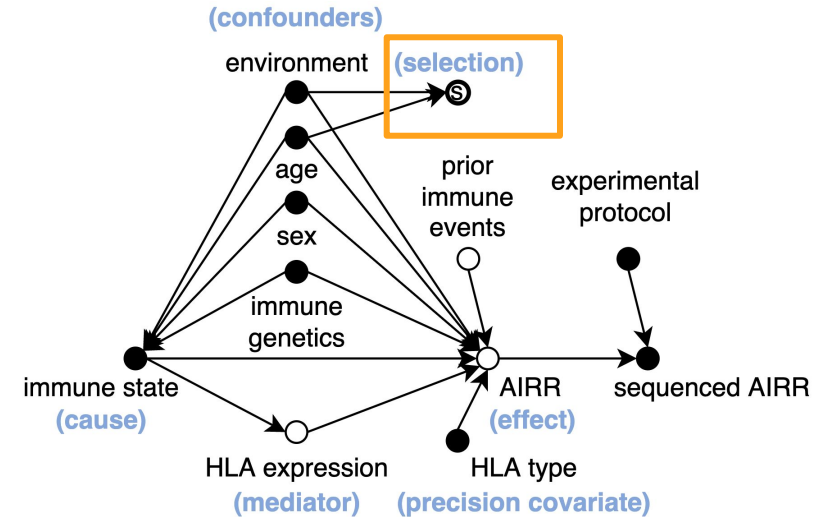
  - ❏ spurious correlations: introduced, removed, reversed



Pavlovic et al. 2022

# The causal inference framework in the AIRR field

- ❏ A causal model for a viral infection (different for different diseases)

- ❏ Selection bias:

  - ❏ preferential selection of study participants

  - ❏ spurious correlations: introduced, removed, reversed

- ❏ Confounding bias:

  - ❏ influence both the immune state and AIRR



Pavlovic et al. 2022

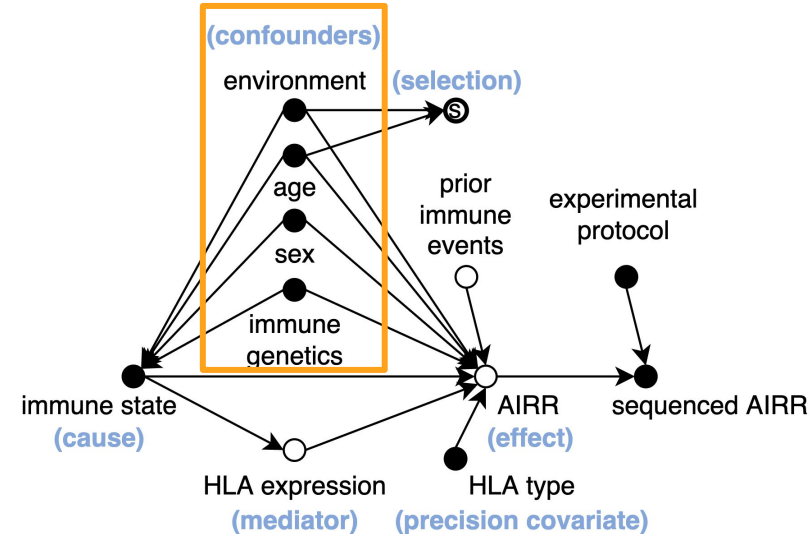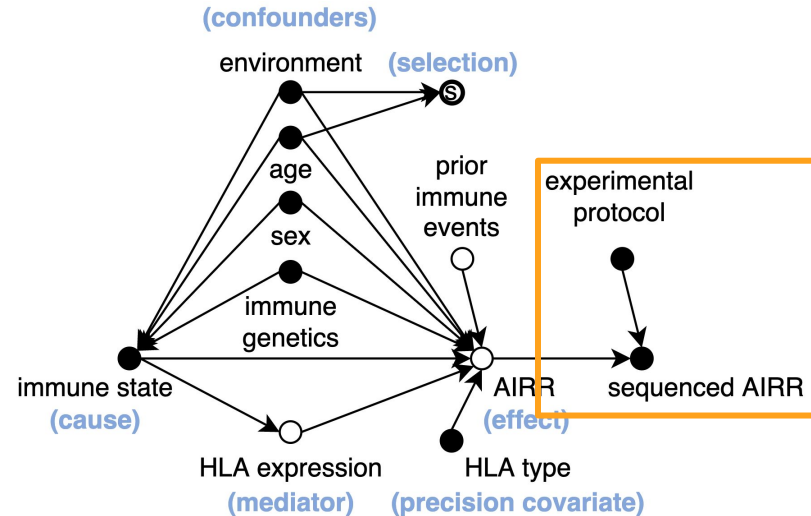# The causal inference framework in the AIRR field

- ❏ A causal model for a viral infection
  (different for different diseases)

- ❏ Selection bias:

  - ❏ preferential selection of study participants

  - ❏ spurious correlations: introduced, removed, reversed

- ❏ Confounding bias:

  - ❏ influence both the immune state and AIRR

- ❏ Batch effects & timing of measurement



Pavlovic et al. 2022

# Summary

❏ AIRR ML ≠ applying several fancy ML method to AIRR data

  ❏ Complex biological structure (both receptors and repertoires)

  ❏ Large variability and sparsity

  ❏ Causal variables (sex, age, HLA etc)

❏ AIRR ML methods should be benchmarked on both experimental and synthetic data with known ground truth

❏ Ultimately we need large-scale experimental data with known ground truth

# Acknowledgements

**UiO Oslo**
Prof. Geir Kjetil Sandve
Prof. Victor Greiff
Prof. Dag T.T. Haug
Prof. Ingrid H. Haff
Prof. Ludvig Sollid
Prof. Torbjørn Rognes
Dr. Fridtjof Lund-Johansen
Dr. Rahmad Akbar
Dr. Igor Snapkov
Dr. Philippe Robert
Dr. Chakravarthi Kanduri
Dr. Ivar Grytten
Dr. Knut Rand
Dr. Enrico Riccardi
Dr. Mai Ha Vu
Dr. Brij Bhusan Mehta
Lonneke Scheffer
Andrei Slabodkin
Ghadi Al Hajj
Robert Frank
Thomas Minotto
Khang Le Quy

**UFlorida**
Keshav Motwani
Prof. Todd Brusko

**ETH Zurich**
Dr. Cédric R. Weber
Prof. Sai T. Reddy

**FHNW**
Prof. Enkelejda Miho

**JKU Linz**
Prof. Günter Klambauer
Prof. Sepp Hochreiter
Michael Widrich

**iReceptor+**
Prof. Gur Yaari
Prof. Lindsay Cowell
Dr. Scott Christley
Dr. Artur Rocha
Alexandre Almeida Costa

**UCSD**
Dr. Yana Safonova
Prof. Pavel Pevzner



**BGI**
Prof. Xiao Liu
Wei Zhang Longlong
Wang Jinghua Wu
Ziyun Wan
Shiyu Wang
Kai Gao