

# TCR distance calculation and clustering: TCRcluster

Zhu Yule, Tang Jiayi, Huang Yutong, Hu Zihao, Lu Jianzhang

Group 1

2022.12.15



## 1 Introduction

## 2 TCR distance calculation

## 3 Clustering

## 4 Visualization

## 5 Verification

## 6 Discussion

# TCR recognition

- Epitope, also called antigenic determinant, is a site or region on an antigen that is recognized by an antibody as a foreign body.
- T cell receptors (TCRs) recognize antigens presented by major histocompatibility complex (MHC) on antigen-presenting cells (APC).

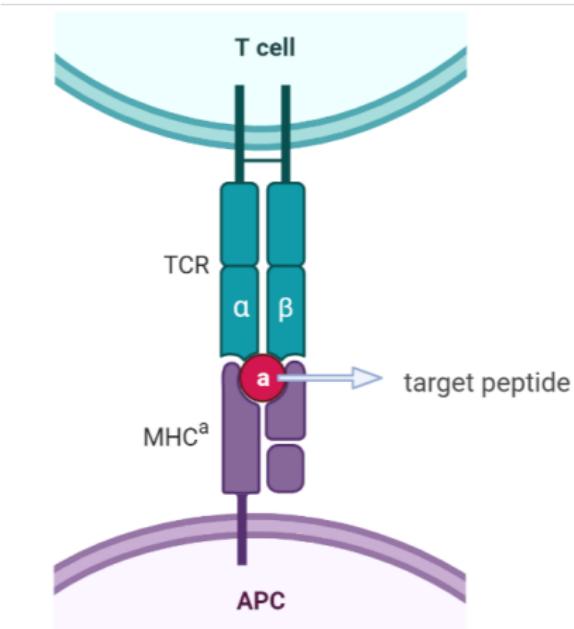


Figure: TCR recognition (created by biorender.com)

# CDR3

- CDR3 regions on  $\beta$  chain may largely decide the recognition of TCR-pMHC
- **Project aim:** Cluster TCR sequences according to their antigen-binding affinity.

The Hypervariable Complementarity Determining Regions (CDRs) of the Antibody Interact with the Antigen

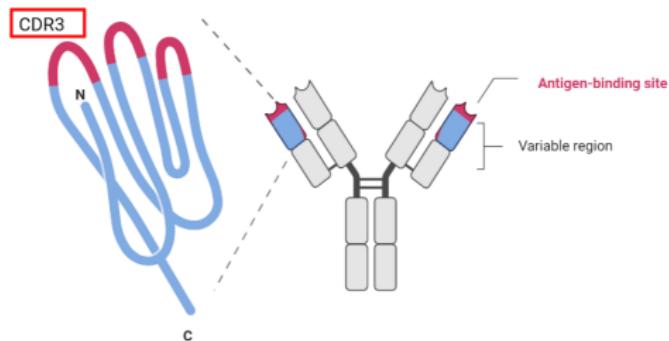
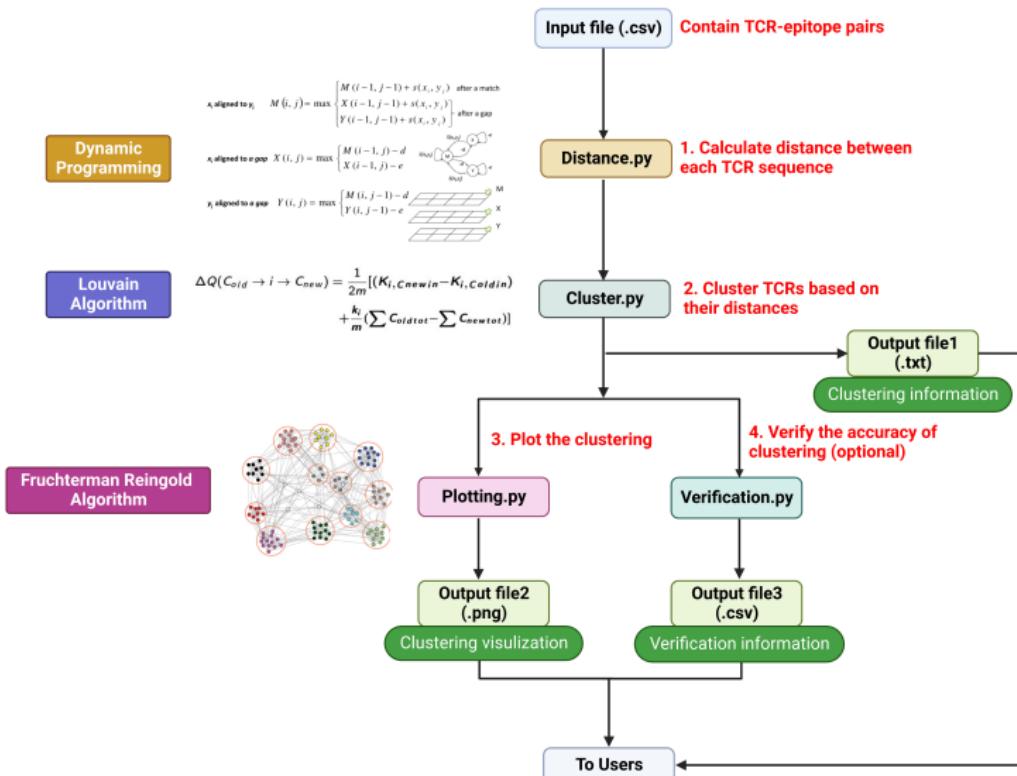


Figure: CDR3 region

# Flow Chat



1 Introduction

2 TCR distance calculation

3 Clustering

4 Visualization

5 Verification

6 Discussion

The similarity (distance) between each TCR sequence

## Workflow



**Figure:** The workflow for calculating distance

## Input

	A	B	C	D	E	F	G
1		CDR3a	CDR3b	HLA	peptide	epitope	MHC
2	seq0	CAAAAGNTGKLIF	CASSRLGASAETLYF	-	HGIRNASFI	-	-
3	seq1	CAAAAYNQGGKLIF	CATSDPAGMTGGWHGYTF	A03:01	KLGGLAQAK	IE-1_CMV	YFAMYQENVAQTVDTLIYIYRDYTWAELAYTWY
4	seq2	CAAADDKIF	CASSQTSIYEQYF	A02:01	LLWNGPMAV	-	YFAMYGEKVAHTHVTDLVRYHYTTWAVLAYTWY
5	seq3	CAAADNYGQNPFV	CAWSSGEGTDTQYF	A11:01	AVFDRKSDAK	EBNA-3B_EBV	YYAMYQENVAQTVDTLIYIYRDYTWAQQAYRWY
6	seq4	CAAETSYDKVIF	CAGGGSQGNLIF	A02:01	GILGFVFTL	Flu-MP_Influenza	YFAMYGEKVAHTHVTDLVRYHYTTWAVLAYTWY
7	seq5	CAAETSYDKVIF	CAGPPVGANNLFF	B08:01	RAFKFQLL	BZLF1_EBV	YDSEYRNIFTNTDESNLYSLYNNYTTWAVDAYTWY
8	seq6	CAAETSYDKVIF	CASSFGNTGLEFF	B08:01	RAFKFQLL	BZLF1_EBV	YDSEYRNIFTNTDESNLYSLYNNYTTWAVDAYTWY
9	seq7	CAAETSYDKVIF	CASSFVDRÄETQYF	A11:01	IVTDFSVIK	EBNA-3B_EBV	YYAMYQENVAQTVDTLIYIYRDYTWAQQAYRWY
10	seq8	CAAETSYDKVIF	CASSFVDRÄETQYF	B08:01	RAFKFQLL	BZLF1_EBV	YDSEYRNIFTNTDESNLYSLYNNYTTWAVDAYTWY
11	seq9	CAAETSYDKVIF	CASSLKGGRHEQYF	A03:01	RLRAEAQVK	EMNA-3A_EBV	YFAMYQENVAQTVDTLIYIYRDYTWAELAYTWY
12	seq10	CAAETSYDKVIF	CASSLOSSOGAPYEQYF	B08:01	RAFKFQLL	BZLF1_EBV	YDSEYRNIFTNTDESNLYSLYNNYTTWAVDAYTWY

**Figure:** The input information

# Dynamic programming (preparation)

- Three states to record the sequences' situation when matching.

## X, Y and M states

X state means to open a gap in the column direction;

Y state means to open a gap in the line direction;

M state means two sequences are matched in this position.

0	-inf	-inf	-inf	-inf	0	-2	-3	-4	-5	0	-inf	-inf	-inf	-inf
-2	0	0	0	0	-inf	0	0	0	0	-inf	0	0	0	0
-3	0	0	0	0	-inf	0	0	0	0	-inf	0	0	0	0
-4	0	0	0	0	-inf	0	0	0	0	-inf	0	0	0	0
-5	0	0	0	0	-inf	0	0	0	0	-inf	0	0	0	0

Figure: X state

Figure: Y state

Figure: M state

# Dynamic programming (formula)

- Dynamic programming for the transformation between states.

## Transformation of three states

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + s(x_i, y_j) & \text{after a match} \\ X(i-1, j-1) + s(x_i, y_j) \\ Y(i-1, j-1) + s(x_i, y_j) \end{cases}$$

AT    TC  
 AT  
 -C  
 -T  
 TC

$$X(i, j) = \max \begin{cases} M(i-1, j) - d \\ X(i-1, j) - e \end{cases}$$

$$Y(i, j) = \max \begin{cases} M(i, j-1) - d \\ Y(i, j-1) - e \end{cases}$$

Figure: The iterative process between states

# Dynamic programming (output)

- Calculating the distance between every two sequences:

The distances for each sequence

```
{Node1:{Node2:dis(1,2),Node3:dis(1,3)...},  
Node2:{Node1:dis(1,2),Node3:dis(2,3)...},  
...}
```

Figure: The output for the sequence distances

The complexity of the algorithm:  $O(n^2 k^2)$

n: the number of the sequences

k: the length of one sequence

1 Introduction

2 TCR distance calculation

3 Clustering

4 Visualization

5 Verification

6 Discussion

# Clustering

- Greedy algorithm  $O(n \log n)$  running time
- According to Louvain algorithm, introduce **Modularity**.
- Modularity  $Q$ : The overall quality of how the network is segmented by communities (**clusters**)

$$Q = \frac{1}{2m} \sum_{ij} [A_{ij} - \frac{k_i k_j}{2m}] \delta(c_i, c_j)$$

$A_{ij}$ : The edge weight between nodes i and j

$k_i$  and  $k_j$ : The sum of the weights of the edges respectively from nodes i and j

$2m$ : The sum of all of the edge weights in the network

$c_i$  and  $c_j$ : The communities of the nodes

$\delta$ : An indicator function  $\delta(c_i, c_j) = 1$  if  $c_i = c_j$  else 0

- $Q \propto \sum_c [(Real\ edges\ within\ c) - (Expected\ edges\ within\ c)]$

# Step 1: Modularity optimization

- Only make local changes to the memberships of nodes-communities.
- At first, each node corresponds to one community.

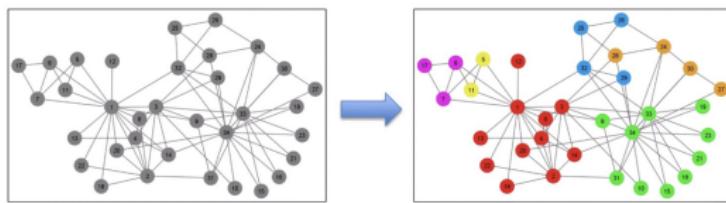


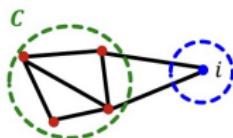
Figure: Modularity optimization step

- Every node → the community of its neighbor.
- A change in the modularity of the partitioning marked as  $\Delta Q$
- $\Delta Q > 0 \rightarrow$  Gain or improve the quality
- $C' = \text{argmax}_c, \Delta Q(C \rightarrow i \rightarrow C')$  Iterate until no node moves.

# Step 1: How to calculate $\Delta Q$ ?

Derive formula:  $\Delta Q(i \rightarrow C)$

Before merging      After merging



$$Q_{\text{before}} = Q(C) + Q(\{i\})$$

$$= \left[ \frac{\Sigma_{in}}{2m} - \left( \frac{\Sigma_{tot}}{2m} \right)^2 \right] + \left[ 0 - \left( \frac{k_i}{2m} \right)^2 \right]$$

$$Q_{\text{after}} = Q(C + \{i\})$$

$$= \frac{\Sigma_{in} + k_{i,in}}{2m} - \left( \frac{\Sigma_{tot} + k_i}{2m} \right)^2$$

$$\Delta Q(i \rightarrow C) = \frac{1}{2m} \left( \frac{\sum_{tot} k_i}{m} - K_{i,i_n} \right)$$

$$\Delta Q(C_{old} \rightarrow i \rightarrow C_{new}) = \frac{1}{2m} [(K_{i,C_{new}} - K_{i,C_{old}}) + \frac{k_i}{m} (\sum C_{oldt} - \sum C_{newt})]$$

## Step 2: Aggregation

The first step → Attain the current stability

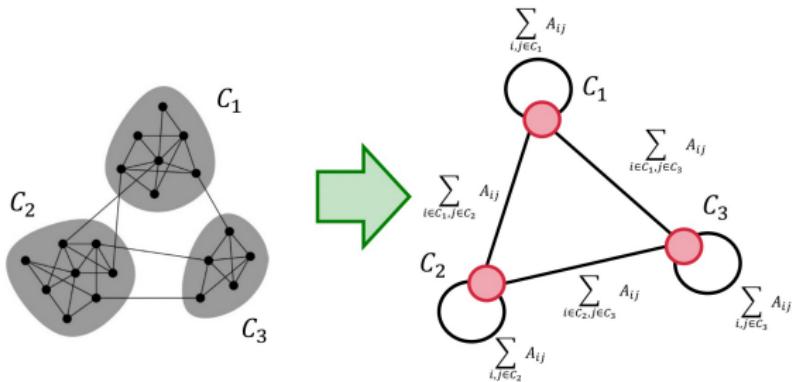


Figure: Aggregation step

- The communities obtained in the first phase are aggregated into super-nodes.
- The edges within the original community → self-loop.
- All edges between the communities → weights between super-nodes.

# Clustering: Summary

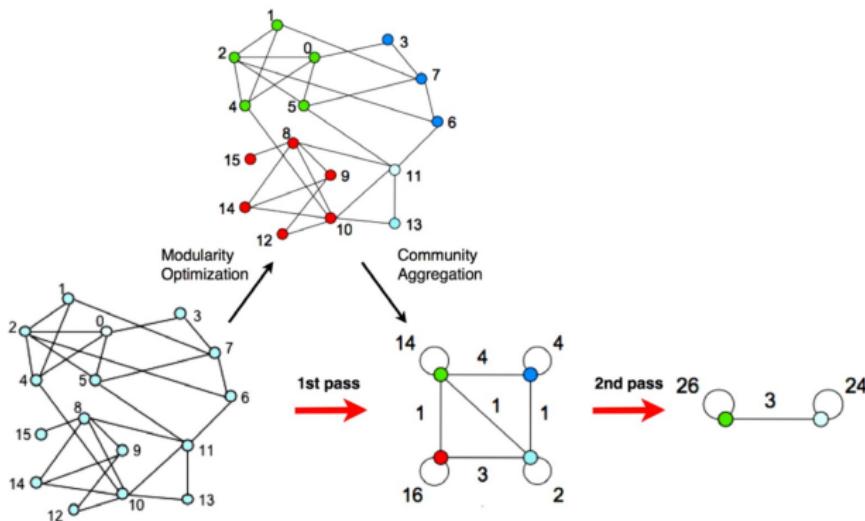


Figure: Summary of the clustering method

- The passes are repeated constantly until no increase of modularity can be gained.

1 Introduction

2 TCR distance calculation

3 Clustering

4 Visualization

5 Verification

6 Discussion

# Visualization: Work Flow

- To preliminarily test the accuracy of our clustering process
- Use **networkx** and **FR algorithm** to do the plotting.

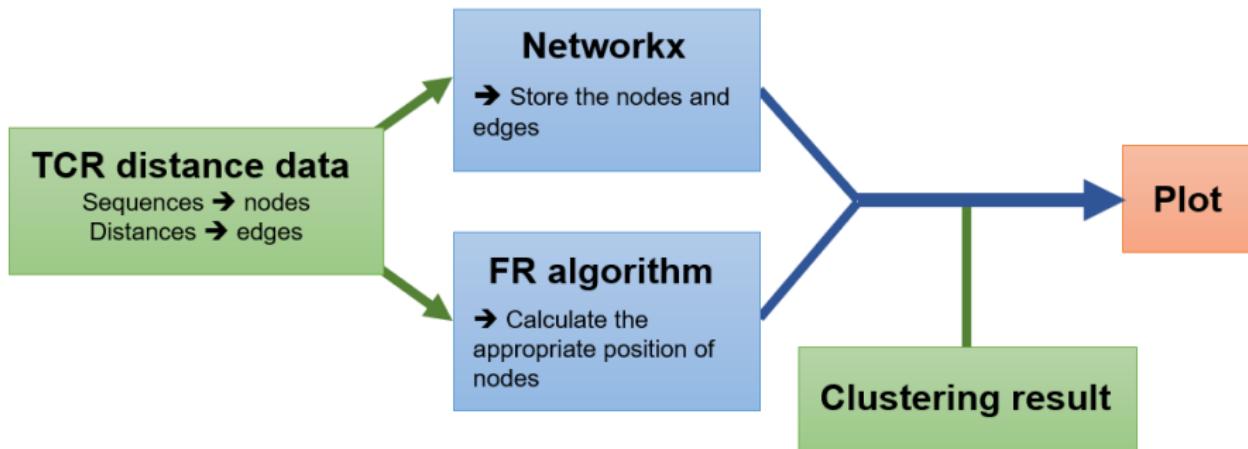


Figure: The visualization work flow

# Visualization: FR algorithm

- **Fruchterman-Reingold algorithm:  $O(n^2)$**

Treat all nodes as charged particles. In this way, there are two forces act on each node: the **coulomb force** and **tensile force**.

- 1) **Coulomb forces:**  $F_q = \frac{k_q \cdot q^2}{r^2}$  keep nodes separated.
- 2) **Tensile forces:**  $F_k = k \cdot x$  promote nodes with high similarity to get closer.

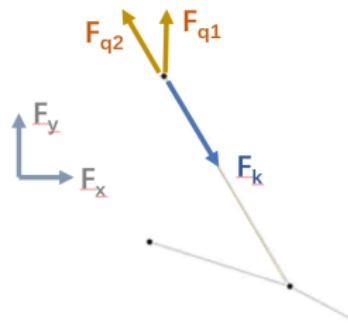


Figure: Force condition analysis in FR algorithm

# Visualization: Result

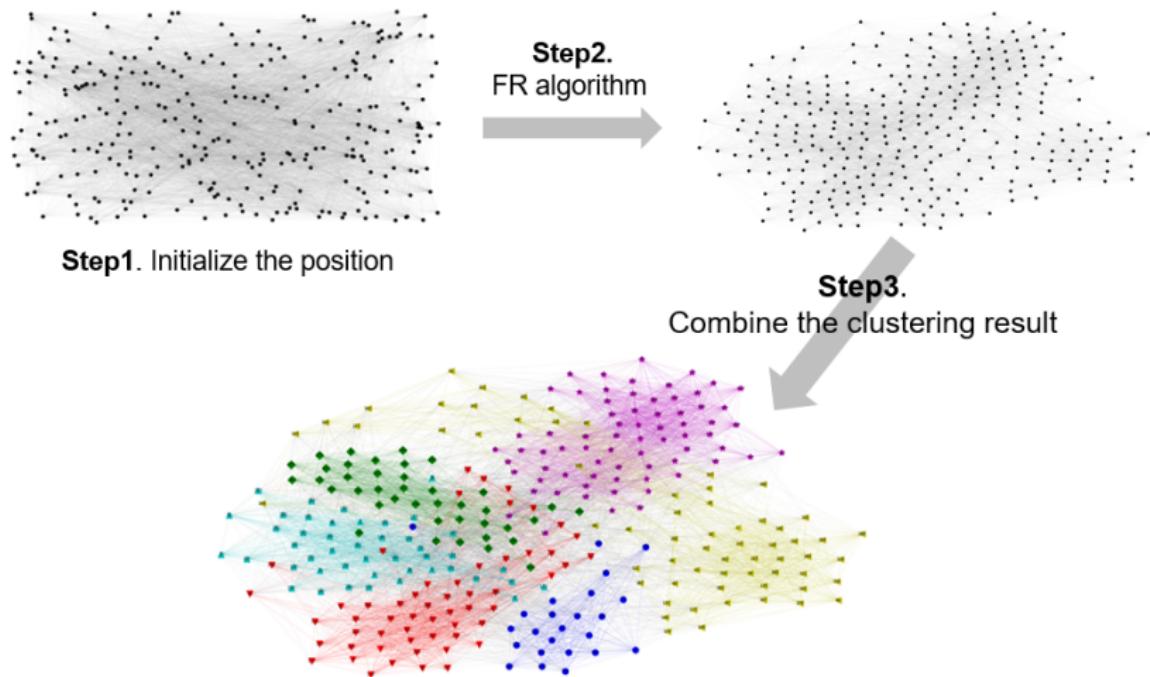


Figure: The visualization process

1 Introduction

2 TCR distance calculation

3 Clustering

4 Visualization

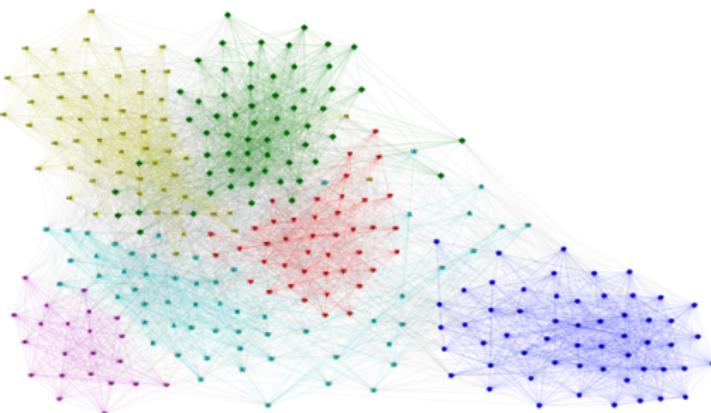
5 Verification

6 Discussion

# Output file

Community 0: 2, 11, 12, 14, 24, 29, 44, 60, 66, 70,  
Community 1: 1, 3, 5, 8, 9, 16, 30, 46, 50, 52, 63, 6  
Community 2: 0, 13, 32, 35, 39, 43, 49, 51, 54, 55,  
Community 3: 19, 22, 25, 27, 33, 61, 65, 71, 74, 90  
Community 4: 15, 20, 28, 40, 41, 42, 45, 48, 59, 86  
Community 5: 4, 6, 7, 10, 17, 18, 21, 23, 26, 31, 34

Output1: clustering information



Output2: visualization

# Verification

- Provide consensus motif for each cluster
- Ideal results:** same or similar antigens in the same cluster

	Community CDR3b	Peptide	Consensus
0	0 CASSIDRGSEAFF	ALSKGVHFV	ALSKGVHFV
1	0 CASSLYGLTEAFF	ALSKGVHFV	ALSKGVHFV
2	0 CASSTEGGTAEFF	ALSKGVHFV	ALSKGVHFV
3	0 CASSRPGGGSTEAFF	ALSKGVHFV	ALSKGVHFV
4	0 CASSLEGQLNTEAFF	ALSKGVHFV	ALSKGVHFV
5	0 CASSYPIWDYTEAFF	ALSKGVHFV	ALSKGVHFV
6	0 CASRGDRGRATEAFF	ALSKGVHFV	ALSKGVHFV
7	0 CAVNRDACKSTF	FLRGRAYGL	ALSKGVHFV
8	0 CAVNRDACKSTF	FLRGRAYGL	ALSKGVHFV
9	0 CASSLGQNTEAFF	FLRGRAYGL	ALSKGVHFV

Output3: consensus motif



Motif logo of the red cluster

## 1 Introduction

## 2 TCR distance calculation

## 3 Clustering

## 4 Visualization

## 5 Verification

## 6 Discussion

# Advantage and limitation of the distance calculation:

## Advantage:

$$\text{Distance} = (D_{CDR3\beta}) * \text{weight} + (D_{CDR3\alpha}) * (1 - \text{weight})$$

CDR3 $\beta$  is required, while CDR3 $\alpha$  also plays a role in antigen recognition.

## Limitation:

Very different (low similarity) TCRs may recognize the same epitope.

## Possible Solution:

Algorithms based on *tcrdist3* or machine learning.

CDR3a	CDR3b	HLA	peptide
seq17027 CAVGGPI	CASTARALNTGELFF	A02:01	LLDFVRFMGV
seq17403 CAVHGYGQNFVF	CVVGEFYF	A02:01	LLDFVRFMGV

Figure: Two TCRs with low alignment score recognize the same peptide

# Advantage and limitation of the clustering:

## Advantage:

1. Do not require setting the number of clusters beforehand.
2. Louvain algorithm greedily maximizes modularity. Our codes run much faster than the standard package.

## Limitation:

1. Unstable result: randomization of the iteration order of the nodes.
2. May yield badly connected communities.

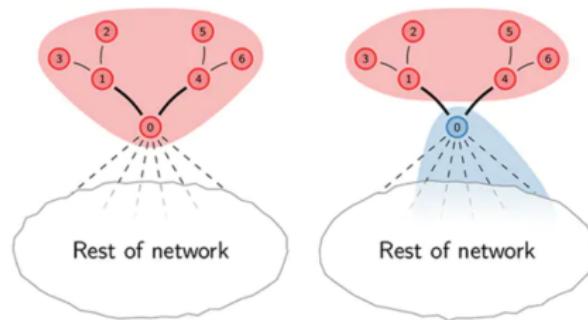


Figure: A reason to yield bad communities

# Advantage and limitation of the plotting:

## Advantage:

FR algorithm is based on the theory of particle physics: easy to understand.

## Limitation:

1. Need to change the position of each node every 0.3 seconds for about 700 times (300 sequences  $\sim 7$  min).
2. Randomization of the initial position of each sequence.

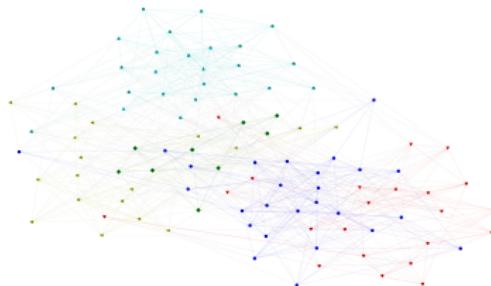


Figure: A bad plotting under 100 sequences

# Significance and application:

## Diagnosis:

1. Reflect the function and immune response state of T cells and design personalized treatments for patients.
2. Detect autoimmune diseases such as type I diabetes.
3. Predict the recognized antigens of new TCRs.

## Novel treatments:

1. Vaccines conferring T-cell immunity against specific cancer mutations: injection of transcribed mRNA coding for antigens.
2. TCR-T cell therapy: Use CRISPR-Cas9 system to knock out endogenous TCR  $\alpha$  and  $\beta$  genes and replace them with tumour-specific TCR sequences.

# References:

-  Blondel, V.D. et al. (2008) Fast unfolding of communities in large networks, *Journal of Statistical Mechanics: Theory and Experiment*, 2008, p. 10008.
-  Dash, P. et al. (2017) Quantifiable predictive features define epitope-specific T cell receptor repertoires, *Nature*, 547(7661), pp. 8993.
-  Foy, S.P. et al. (2022) Non-viral precision T cell receptor replacement for personalized cell therapy, *Nature [Preprint]*.
-  Fruchterman, T.M.J. and Reingold, E.M. (1991) Graph drawing by force-directed placement, *Software: Practice and Experience*, 21(11), pp. 11291164.
-  Glanville, J. et al. (2017) Identifying specificity groups in the T cell receptor repertoire, *Nature*, 547(7661), pp. 9498.
-  Traag, V.A., Waltman, L. and Van Eck, N.J. (2019) From Louvain to Leiden: guaranteeing well-connected communities, *Scientific Reports*, 9, p. e5233.
-  Tusup, M. et al. (2021) mRNA-Based Anti-TCR CDR3 Tumour Vaccine for T-Cell Lymphoma, *Pharmaceutics*, 13(7), p. 1040.
-  Stanford. Stanford cs224w: Community detection in networks. 2021. (Group 1)

*Thanks for listening*