# Abalign   Manual
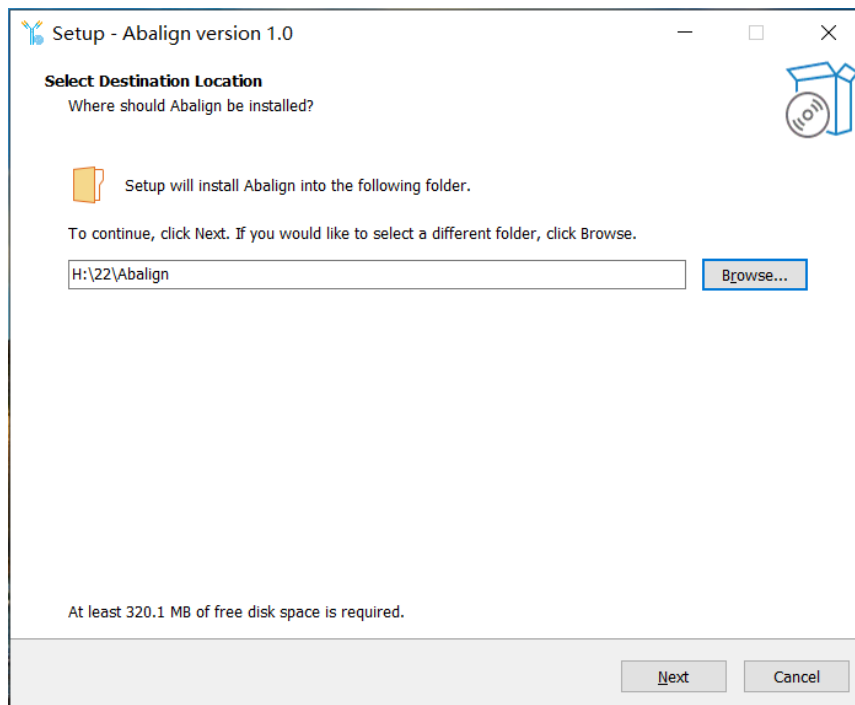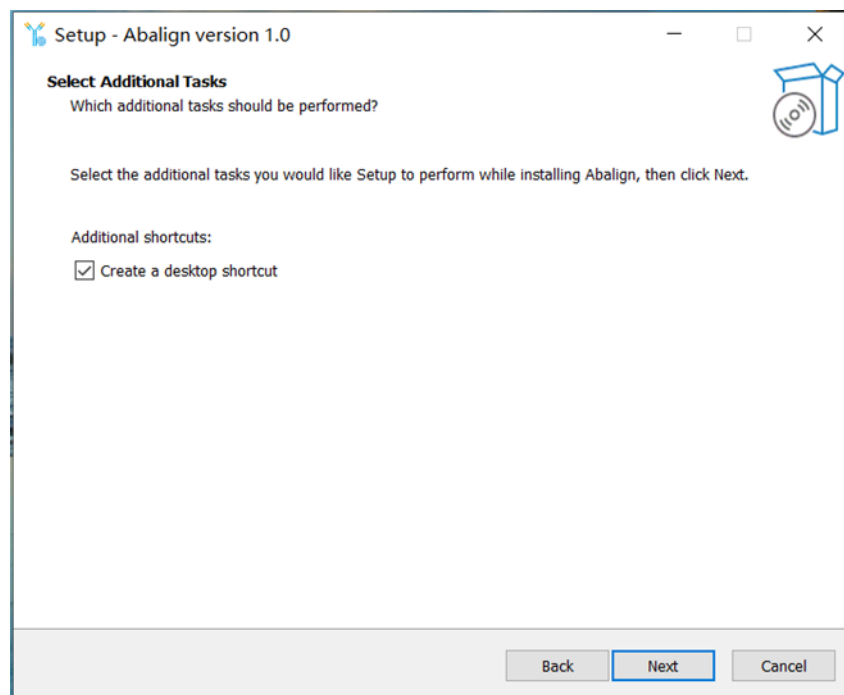
-

-

# Introduce

Multiple sequence alignment has long been used as a powerful tool to investigate the evolutionary, structural and functional properties of protein families. It is also a fundamental technique in recent deep-learning based protein 3D structure prediction methods. The existing multiple sequence alignment methods are extremely difficult to align the highly variable regions of antibody or B-Cell Receptor (BCR) sequences without the prior knowledge of antibody gene recombination and hypermutation in the process of maturation. This multiple sequence alignment tool, named Abalign, which integrates heuristic knowledge of the standardized antibody sequence numbering schemes, including IMGT, KABAT and Chothia systems. The alignment follows the well-characterized patterns of the conserved or highly variable positions known by immunology studies. Hence the alignment result is consistent with the structural and immunological knowledge. Abalign was implemented in a user-friendly software with interactive and visual interface, which supports the multiple sequence alignment, as well as sequence clustering, antibody numbering, CDR delimiting, constructing phylogenetic tree, sequence similarity heat map, V-gene determination and abundance analysis by just clicking the buttons. Abalign allows the high-throughput analysis for BCR sequencing data, which can be finished in 850 minutes for 1 Gb DNA fasta sequences through a single thread by a notebook equipped with I7 10870H. Abalign is a powerful and efficient tool for biological researchers to analyze massive amounts of BCR or antibody sequences and get new discovery in immunoinformatic study.
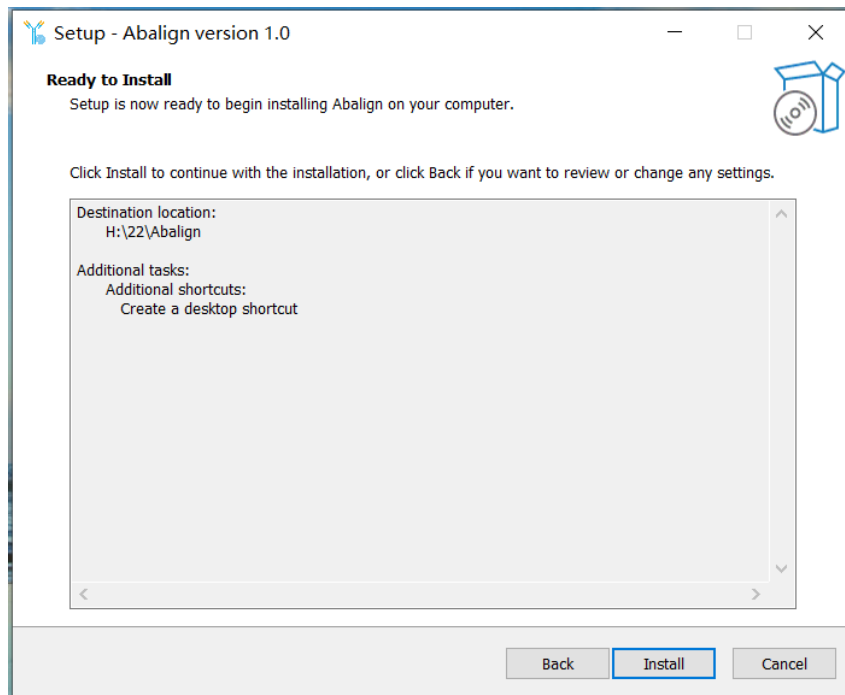
# Setup (for windows)

Double click **Align Setup.exe** to enter the installation, click **Browse** to select the installation path, and click **Next**.

Select whether to **Create a desktop shortcut**, click **Next**, and click **Install** to install.

Please wait for the installation to complete, select whether to **Launch Abalign**, and click **Finish** to exit the installation program.



If you need to uninstall the software, click **unins000.exe** under the installation folder or uninstall or change the program through windows to complete the uninstall.

# Usage

## Software constitution

**Note: Do not modify any content of lib folder**

**Executable file (UI)**: The file named Abalign_ui in the software root directory is the visual window program, double-click to run.

**lib folder**: Do not make any changes to the lib folder under the software directory, which contains a series of configuration files and binary files that implement multiple sequence alignment.

**example folder**: The sequence files of DNA and amino acids in FASTA format are stored for user testing.

## Basic functions

**File input:** The program supports the amino acid sequence file or nucleic acid sequence file in FASTA format. The input file can be selected by clicking the **Input** button or dragged directly into the text box to input the file.

**Numbering strategy:** Select the antibody numbering strategy (support Chothia, Kabat, IMGT strategy) and the type of antibody sequence by clicking the "**IMGT for heavy chain**" drop-down menu (select the heavy chain will detect all the antibody heavy chains in the input sequence and filter out the light chain and the non-variable region sequence, the light chain is the same)

**Multiple sequence alignment:** Click the **Align** button to find the variable domain of the antibody and use sequences to execute multiple sequence alignment. During the alignment, the V Gene and the species that are most similar to each variable region sequence are found, and the V Gene name and species name are displayed after the corresponding sequence name. After running, the residue number of each site will be displayed above, and the seven regions of the antibody variable domain (FR1, CDR1... CDR3, FR4) are rendered into different colors by default. If you need to change the rendering method, please click **Display- > Sequence render mode**.
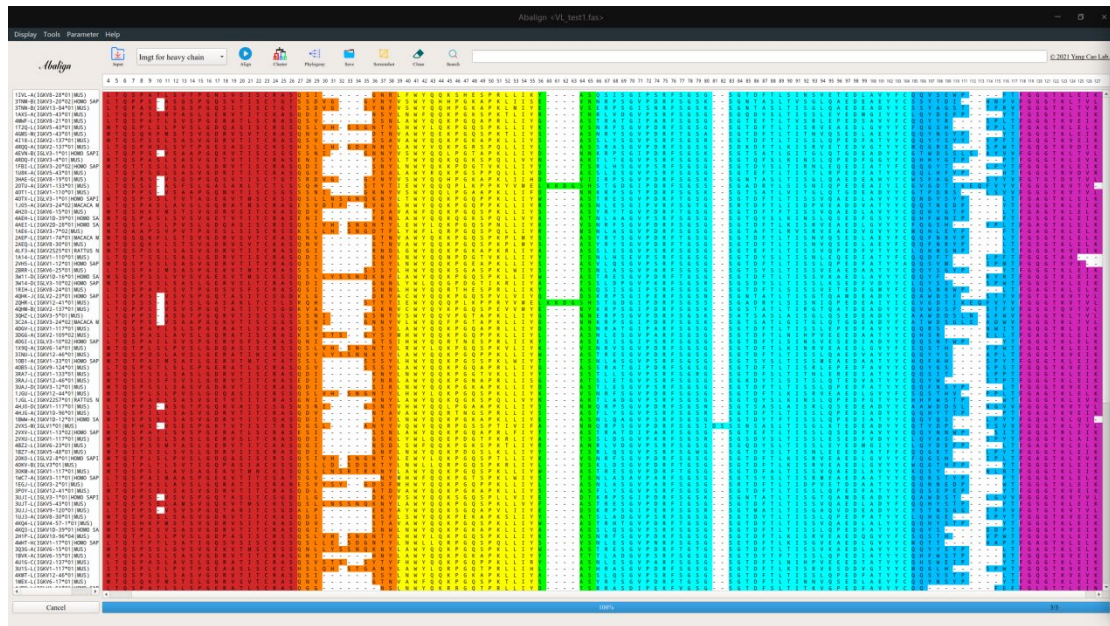
Figure 1. Multiple sequence alignment results.

**File Save:** Click the **Save** button to save the contents of the current interface as .fasta file and .temp file. The former stored the multiple sequence alignment results and the latter stored the multiple sequence alignment results with '*' as a separation of the various regions in the variable region.

**Terminate process:** Click the **Cancel** button to terminate the ongoing multiple sequence alignment.

## Optional functions

**Note: All optional functions require multi-sequence alignment first.**

**Hierarchical clustering:** Click the **Cluster** button to run hierarchical clustering. This function will hierarchically cluster the multiple sequence alignment data after Align. After the clustering, the hierarchical clustering tree will be displayed. The sequence will be rearranged according to the tree, and the same sequence will be arranged adjacently, and the sequence name will be rendered as the same color. The parameters of hierarchical clustering can be changed by the **Parameter** option in the menu bar.
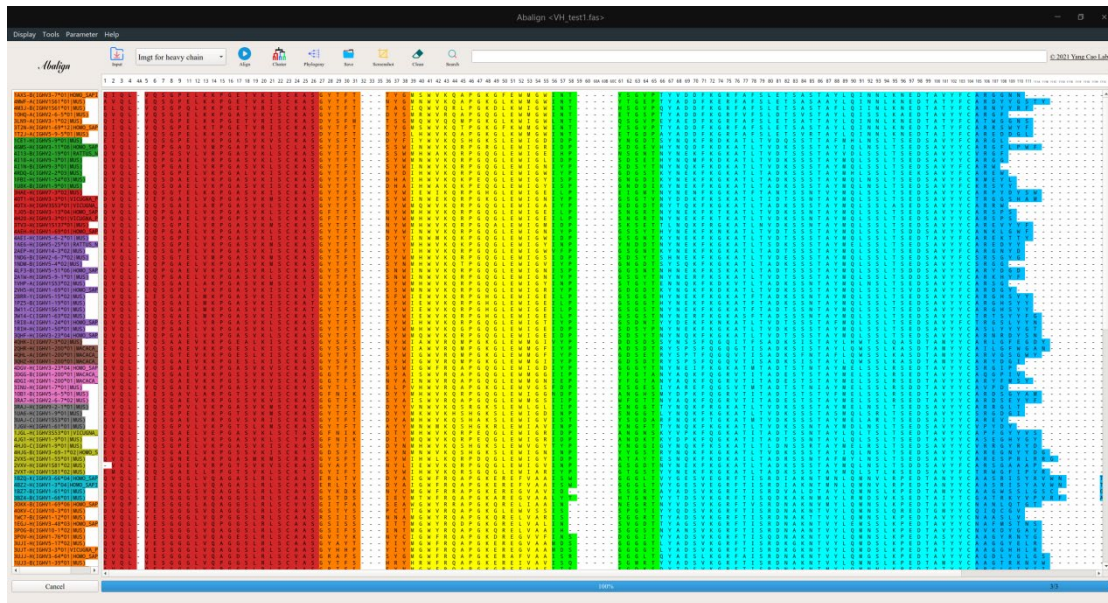
Fig. 2. Hierarchical clustering results of antibody variable domain

**Phylogenetic tree:** This function uses FastTree software ( maximum likelihood method )[1] to build .nwk file and visualize it with Ete3. After clicking **Phylogeny**, the parameter box of the phylogenetic tree will be popped up. After clicking **Run**, the phylogenetic tree will be constructed. The text box below will enter the log information. Visualization results of the phylogenetic tree appear after running. If you need to save the .nwk file of the current phylogenetic tree, click the **Save** button below the parameter box to save.
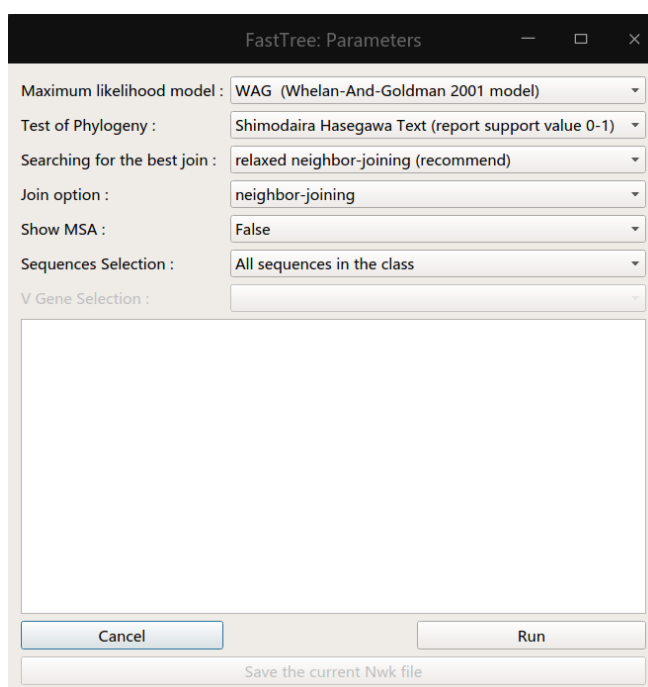


Fig. 3. Parameter panel of phylogenetic tree

(a) Maximum likelihood model : Maximum likelihood model for amino acid substitutions, including WAG, JJT, LG

(b) Test of Phylogeny : branch checking methods for phylogenetic trees, including Shiodaira Hasegawa Text(defined by FastTree) and Bootstrap(1000x)

(c) Searching for best join : Using the neighbor-join method to construct a rough tree topology of the system, the default relaxed neighbor joining is faster, the exhaustive search is slower, and the search of the visible

set is only the fastest

(d)　Join option : The neighbor-join method used to construct a rough topology, defaulting to neighbor - join, in addition to BioNJ.

(e)　Show msa : the results of multi-sequence alignment will be displayed with the evolution tree.

(f)　Sequences Selection: By default, All sequences in the class, the system tree is constructed with all sequences currently displayed, Sequences belong to a specific V Gene, and the system tree is constructed with sequences belonging to V Gene selected in V Gene Selection.

(g)　V Gene Selection: The option contain the V Gene types of all sequences on the current page. Availabled When option of "Sequences Selection" switch to "Sequences belong to a specific V Gene",
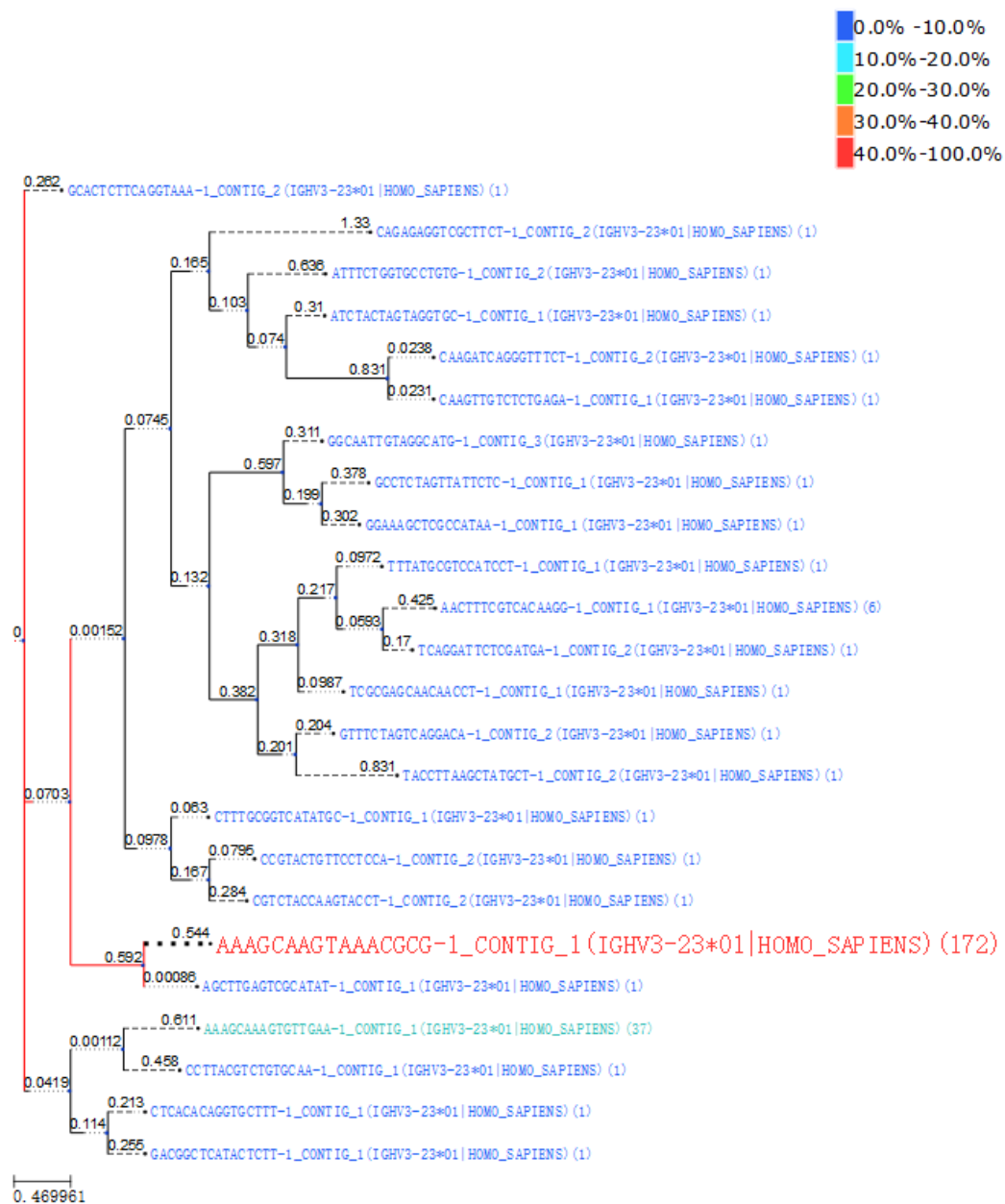


Figure 4. Visualization results of system tree constructed according to V Gene
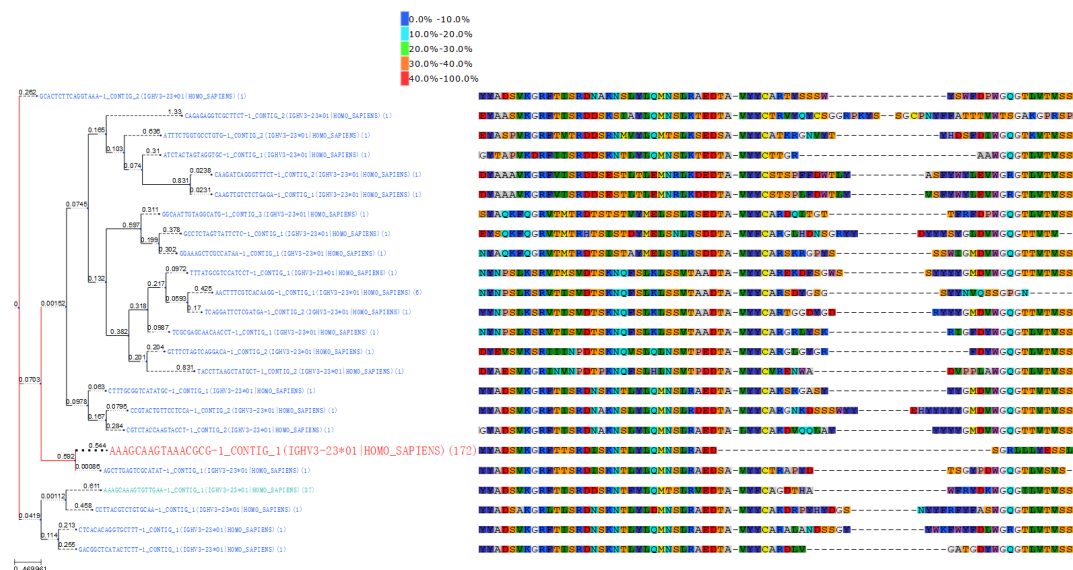
Fig 5.　Phylogenetic tree with multiple sequence alignment results

In this program, the root node is automatically set to the sequence most similar to the germline V Gene in all sequences of the phylogenetic tree, and the germline V Gene is the V Gene of the most abundant sequence. The red branch in the tree represents the evolutionary path from the root node to the sequence with the highest abundance. The program renders the label of leaf nodes by the abundance of sequences. The higher the proportion of the number of sequences of leaf nodes in the total number of sequences of all leaf nodes is, the closer the color of leaf node label is to red and the larger the font is.

**Interface screenshot:** Click the **Screenshot** button to save the result screenshot, which is in .png format by default. Note: Too many sequences may not be successful.

**Character Search:** Enter the string you want to query in the text box after the **Search** button. Click the **Search** button to search for the location of the string in the sequence name and sequence.
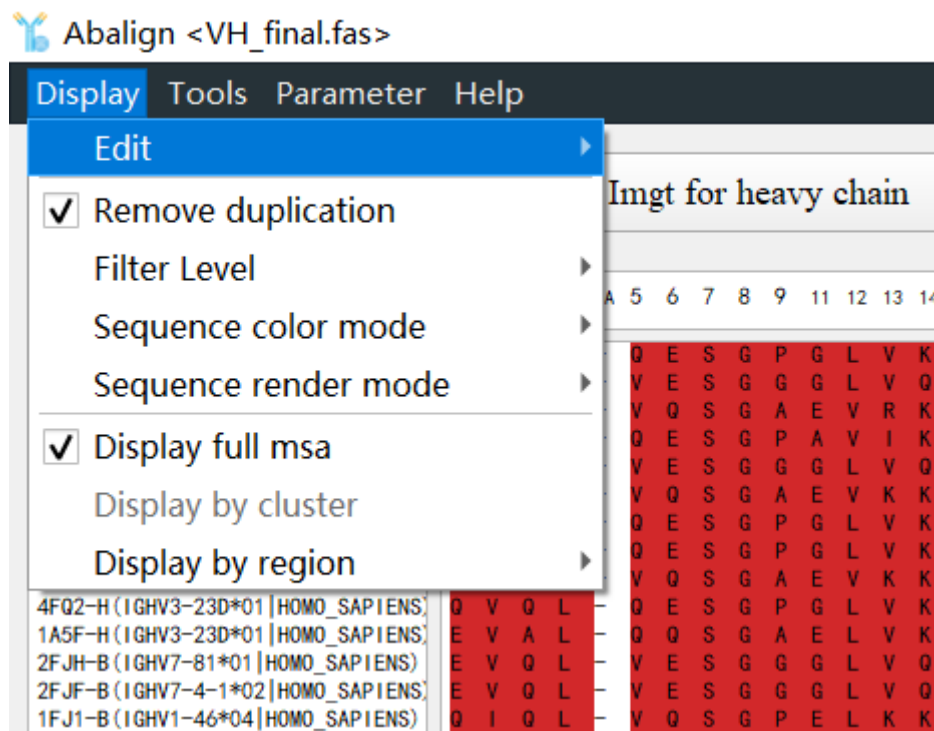
## Menu bar

### Display



Fig 6.   Display options

**Remove duplication:** By default, repeats in the variable region sequence of the antibody are removed during Align.

**Sequence filtering:** "**Normal**" is defaulted. After selecting "**Normal**" or "**Strict**", the length of each region in the variable region of the antibody will be limited in the Align process. After selecting "**Off**", the restriction will be closed. If it is not within the limit, the sequence will be filtered out. "Strict" is more restrictive than "Normal".

**Switch sequence rendering mode**: Click "**Display**" and move the mouse over "**Sequence color mode**", "**Light mode**" (more vivid color) and "**soft mode**" (more soft color) can be selected. Click "**Display**", move the mouse to "**Sequence render mode**", and select "**Color by region**" (different regions of the antibody sequence rendered as different colors) and "**Color by amino**" (different amino acids rendered as different color).

**Select the display area:** click "**Display**", select Display full msa (default) to display the entire msa results; move the mouse to the "**Display by region**", you can choose to display different regions of the antibody; click "**Display by cluster**" (after complete the cluster operation)

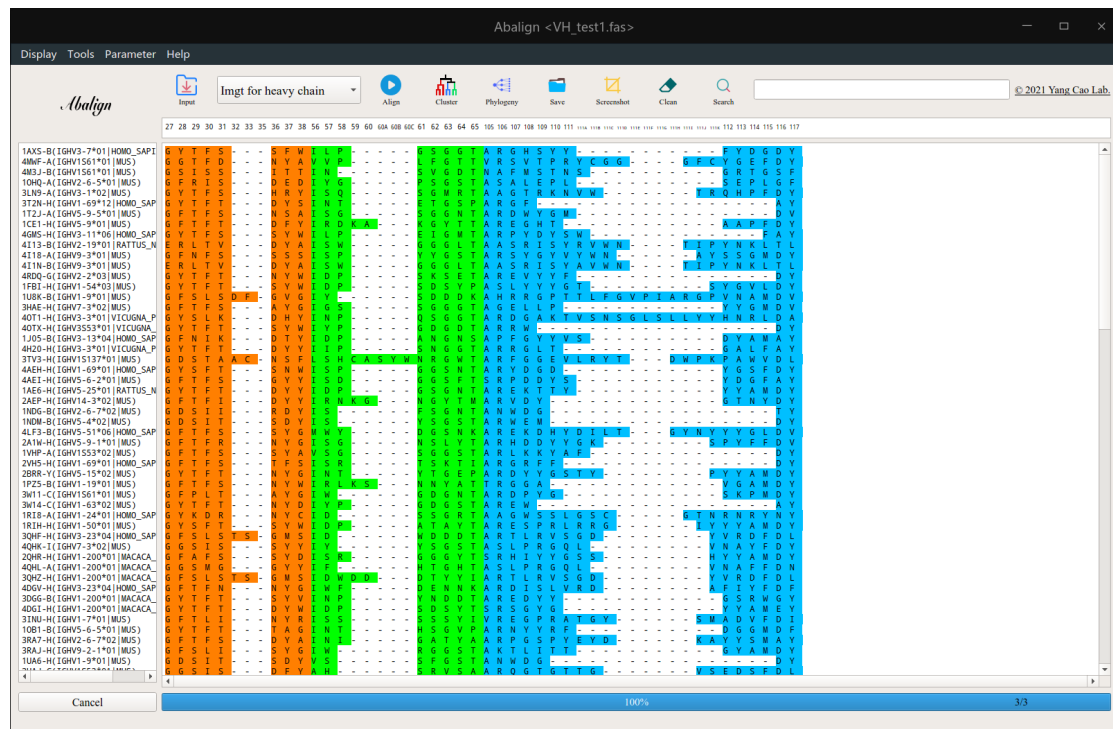to select the class that needs to be displayed by pop-up the tab.



Figure 7. Display only the CDR1, CDR2, CDR3 regions of the antibody variable domain

## Tools

**Find V Gene:** The default selection cannot be changed to find the most similar V Gene for each sequence of antibody variable domain and the species to which the V Gene belongs, and display the V Gene name and the species name behind the corresponding sequence name.

**Heatmap:** Click on "**Heatmap**" to generate a heatmap based on the similarity of the current sequence, and you can cluster the heatmaps and wait for a while if the number of sequences is large. All Heatmaps are based on the content of the current multi-sequence alignment text box. If "Display" is used to change the current displayed sequence, Heatmap will change accordingly.
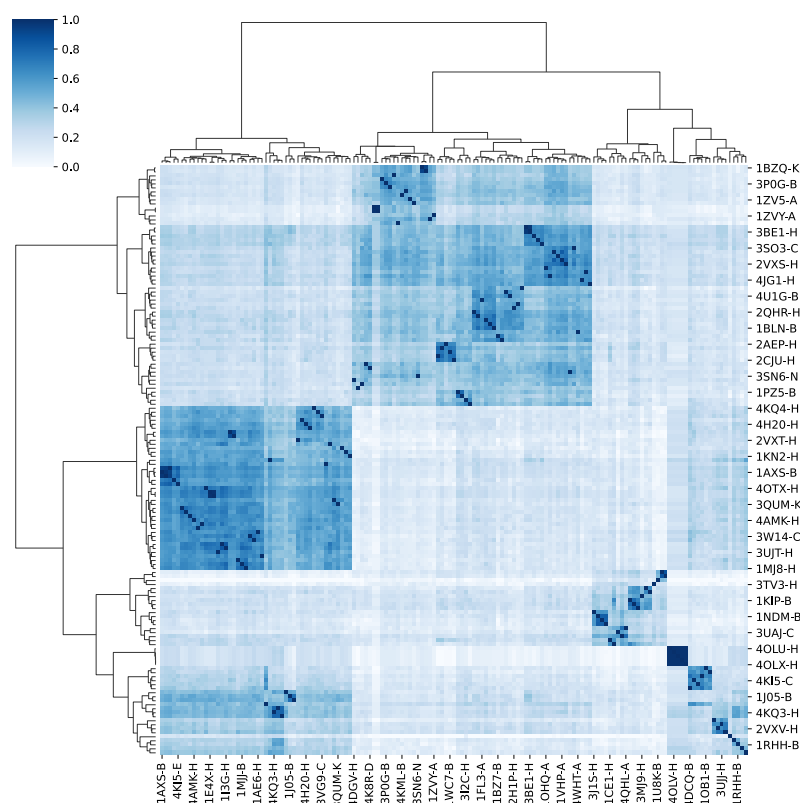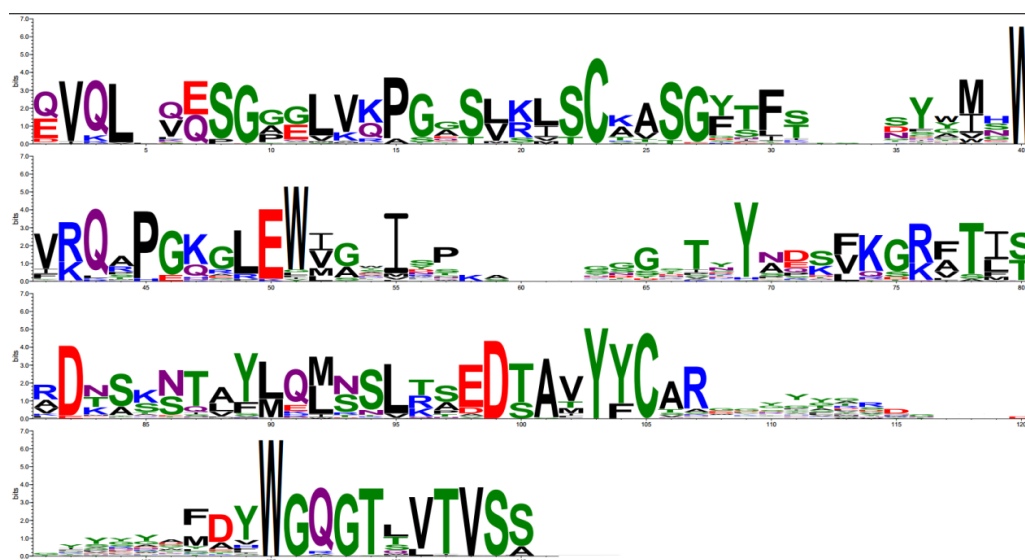
Figure 8. Heatmap after hierarchical clustering

**Seqlogo:** Click **Tools**, move the mouse to "**Seqlogo**", click "**By Entropy**" to get entropy ordinate Seqlogo, click "**By Frequence**" to get frequency ordinate Seqlogo. Move the mouse to the "**Color**" option to change the rendering mode of Seqlogo. All Seqlogo is based on the content of the current multiple sequence alignment text box. If Display is used to change the current displayed sequence, Seqlogo will change accordingly. To ensure compatibility, all Seqlogo is in



PDF format.

Figure 9. Seqlogo for the entire antibody variable region, with entropy as the y-axis

**Abundance map:** Click "**Tools**", move the mouse to "**Abundance**", see "**V Gene abundance**" and "**Sequence abundance**" options, and generate V Gene abundance and Sequence abundance maps after clicking.
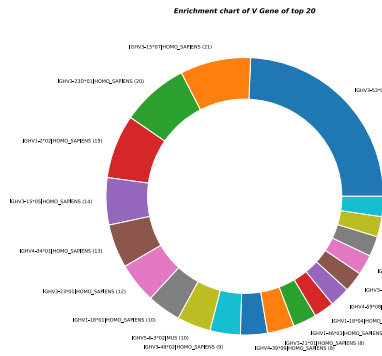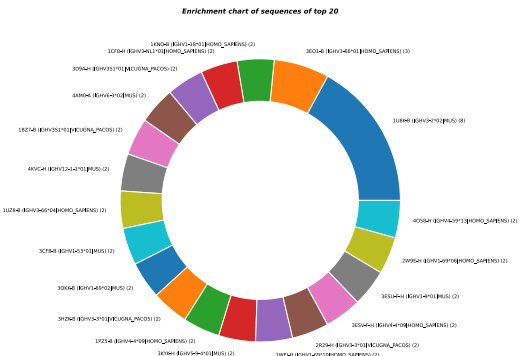


Figure 10. V Gene abundance map



Figure 11. Sequence abundance map

## Parameter

After clicking the "**Parameter**" in the menu bar, a dialog box will pop up, which can modify the clustering penalty, clustering weight ( if the weight of CDR3 region is increased, the similarity of CDR3 region will be paid more attention to in the clustering process ) and the parameters of Hierarchical Clustering Dendrogram.



Figure 12. Hierarchical Clustering Dendrogram parameters

(a) Branch: It represents the number of branches displayed in hierarchical clustering tree. This parameter takes effect with truncate_mode parameter.

(b) Threshold: represents the threshold of clustering partition. If this value is 0.1, a tenth of the total length of the Y-axis (the maximum difference between sequences) is used as a threshold to partition different classes. The more the number of clusters, the smaller the value should be set.

(c) leaf _ font: Represents the size of the branch label font

(d) cluster_method: represents the method used to calculate the distance between two classes during clustering.

(e) root_orientation: Represents the direction of the root.

(f) count_sort: indicates whether the clustering results are sorted according to the number of each category.
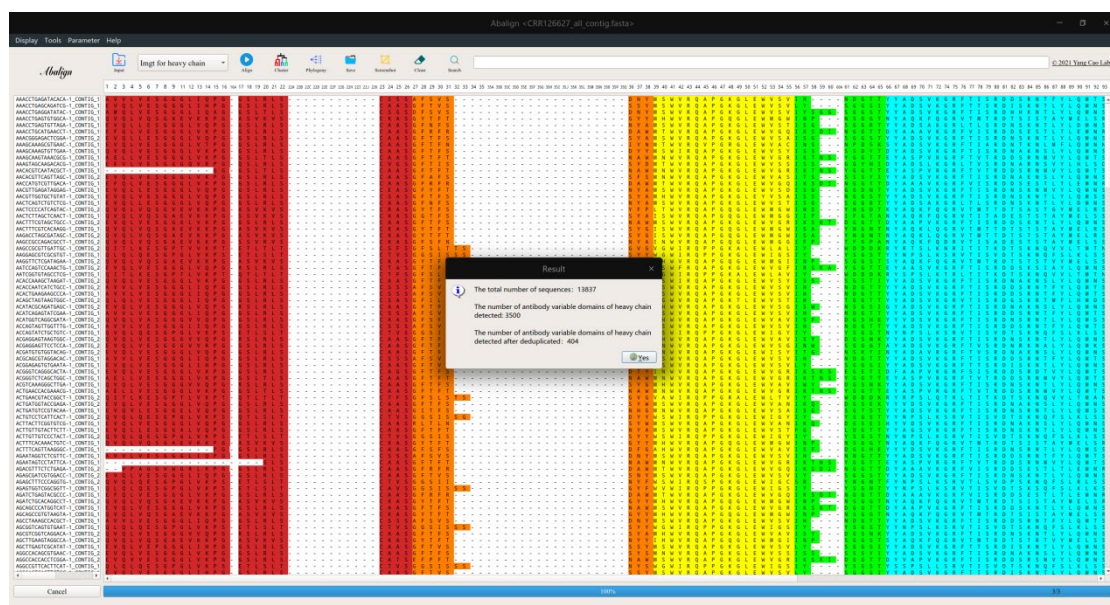
(g) Distance_sort: indicates whether the clustering results are sorted according to the distance of each category.
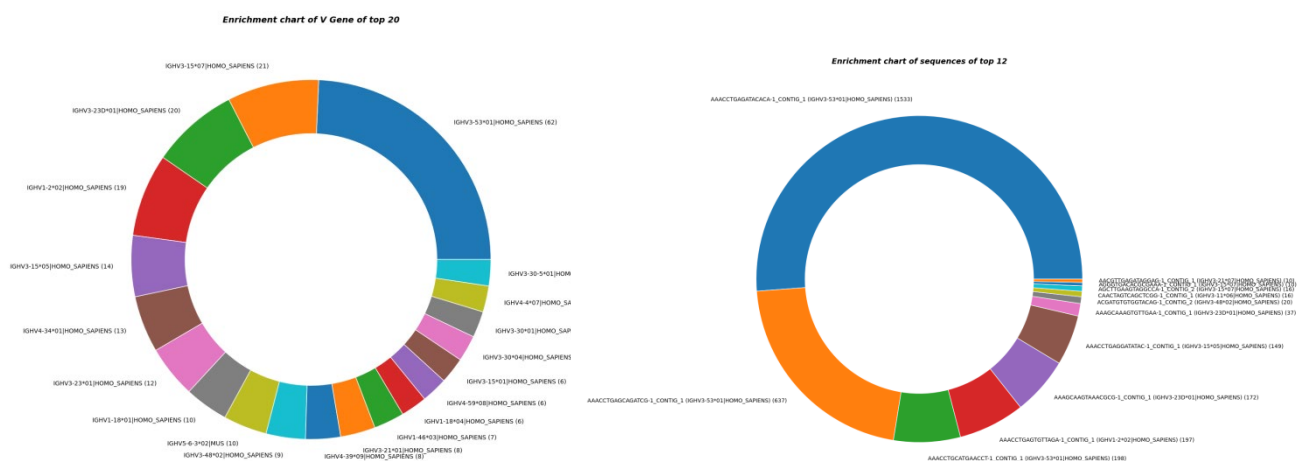
# Example

In this case, we used the samples obtained from single cell sequencing of patients with new coronavirus for demonstration.

**Step 1 Input file:** select the fasta file that needs to be processed after clicking Input.

**Step 2 Search for antibody variable domain and multiple sequence alignment:** After file loading is completed, click "Align" to run the program, the progress bar below the program will show the progress of the program, if the file is too large, the progress bar is not updated in a short time is normal. After the comparison, a dialog box will be popped, and the total number of input sequences will be displayed in the dialog box. The number of sequences in the variable region of the antibody and the number of sequences after deduplication will be detected.



**Step 3 looks at multiple sequence alignments:** click the "Tools" button in the top menu bar and

move the mouse to "Abundance". Click "V Gene abundance" or "Sequence abundance" to obtain the abundance map of V Gene and sequence. Somatic hypermutation occur during antibody maturation and the final mature antibody is highly expressed in the body. Therefore, it is meaningful to study antibodies with high abundance.

**Step 4 Build the phylogenetic tree:** according to the sequence abundance information, build the tree with the V Gene of the high abundance sequence. In this case, the most abundant sequence belongs to V Gene "IGVH3-53 * 01". Click Phylogeny to open the parameter list of tree building, switch the Sequences Selection option to "Sequences belonging to a specific V Gene", select "IGVH3-53 * 01" in the V Gene Selection option, and click Run to start building the evolutionary tree. If you are satisfied with the tree building results, click the "Save the current Nwk file" button in the tree building menu to save the .nwk file of the current system tree.
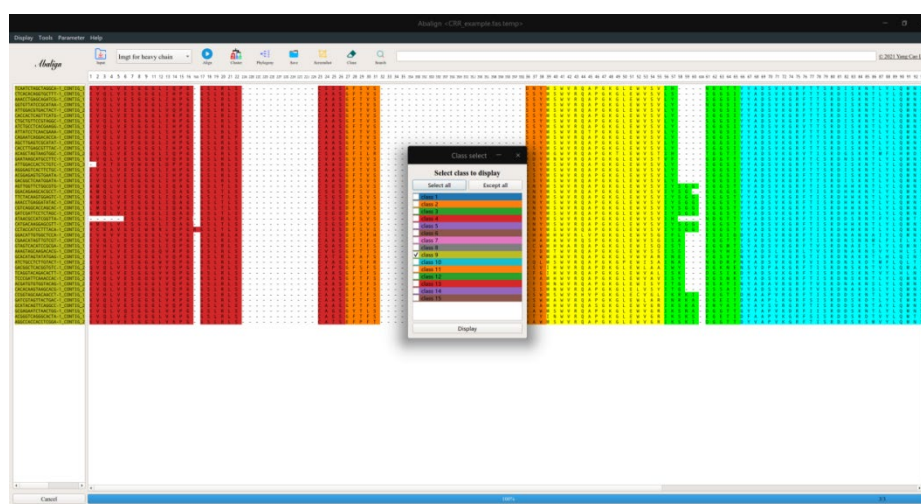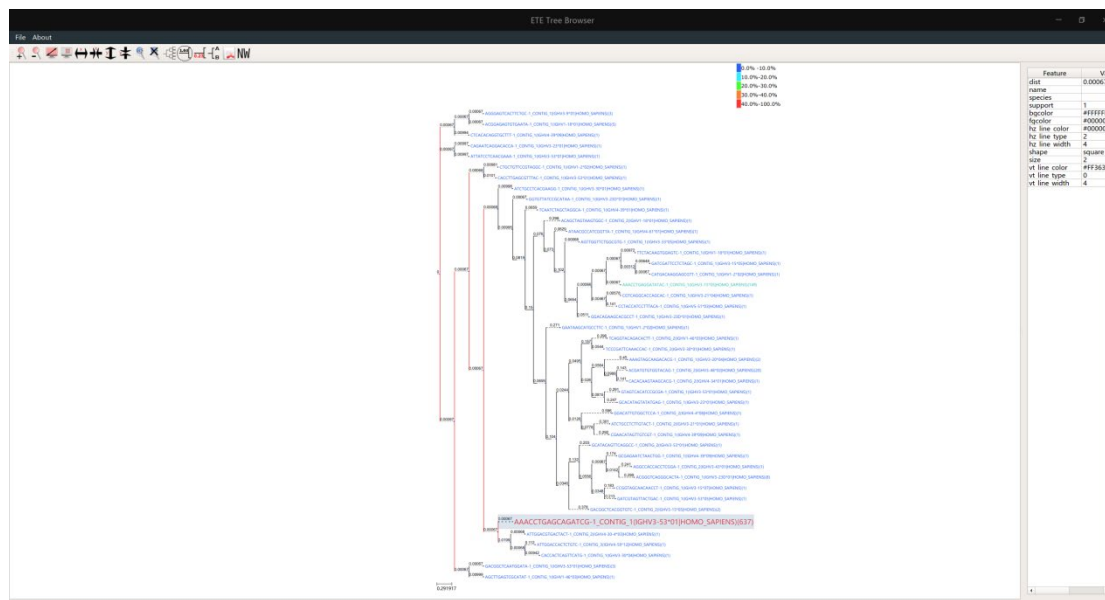
**Step 5 View the phylogenetic tree:** Once the evolution tree is built, a visual window pops up that adjusts the view and displays other information about the tree through the button above the visual window. After selecting the node of the tree, you can also modify the attributes of the node, such as the color of the node, in the right window.



**Step 6 Hierarchical clustering:** After clicking Cluster, hierarchical clustering can be carried out. If the file is too large, it will consume more time, please wait patiently. After the clustering, the hierarchical clustering tree graph will be displayed. Multiple sequence alignment will be rearranged according to the results of the tree graph, and the sequence name will be rendered as the same color according to the color of the branch of the tree graph. After the clustering is completed, click Display- > Display by cluster in the menu bar to select the class that needs to be displayed.

**Step 7 Building phylogenetic trees with different classes:** After selecting one or more classes that need to be used to build trees in Display by cluster, click Phylogeny again, adjust "Sequences Selection" to "All sequences in the class", and click Run to build phylogenetic trees.



# Notice

1. If there is no response during the operation of the program, wait a little and the program is still running.

2. Multiple sequence alignment is for input files. If multiple sequence alignment is carried out, then multiple sequence alignment is carried out again, or multiple sequence alignment is carried out for input files, rather than after alignment.

3. Clicking the Cancel button can only terminate the alignment process and cannot be canceled by Cancel when the alignment sequence is finally loaded.

4. Hierarchical clustering consumes resources and time. To hierarchically cluster a large number of sequences, ensure that the computer has sufficient memory.

5. If the software screen is not displayed properly, adjust the scale of the screen.

This software is developed by Yang Cao Laboratory, College of Life Sciences, Sichuan University. The main developers are Fanjie Zong, Chenyu Long, Wanxin Hu, and Yang Cao. If you have any opinions or suggestions, please contact cy_scu@yeah.net.

# Reference

[1]  Price MN, Dehal PS, Arkin AP. FastTree: computing large minimum evolution trees with profiles instead of a distance matrix[J]. Mol Biol Evol. 2009 Jul;26(7):1641-50.