**M EPITOPIC**

---

**Epitope Fingerprinting Report**

Date: October 23rd, 2024

Work Package:  Epitope Fingerprinting for one monoclonal antibody

Task: Epitope identification for anti-mouse CD184 (BD551966, Rat IgG2b)

---

### General comments

The CD184 antibody was subjected to epitope fingerprinting. The epitope was determined from using NGS data of peptide phage of the naïve ENTE-1 library selected on the antibody coupled to M270 Carboxy magnetic beads (EDC/NHS coupling).

A sequence pool of 142,381 resp. 244,407 sequences from a first and a second selection round was analysed for statistical enrichment of motifs corresponding to potential epitopes in the protein. An excellent enrichment of sequences corresponding to a discontinuous epitope was found.

This report is based on data analysed earlier. In the last years epitopic has improved many statistical analysis tools. The analysis of the data is now not only faster, but we are also able to produce even more details of the antibodies' binding to the antigen.

In an internal challenge we managed to run the analysis with more details revealed in a much shorter time. With the excellent datasets it took only about eight hours for the entire work presented here.
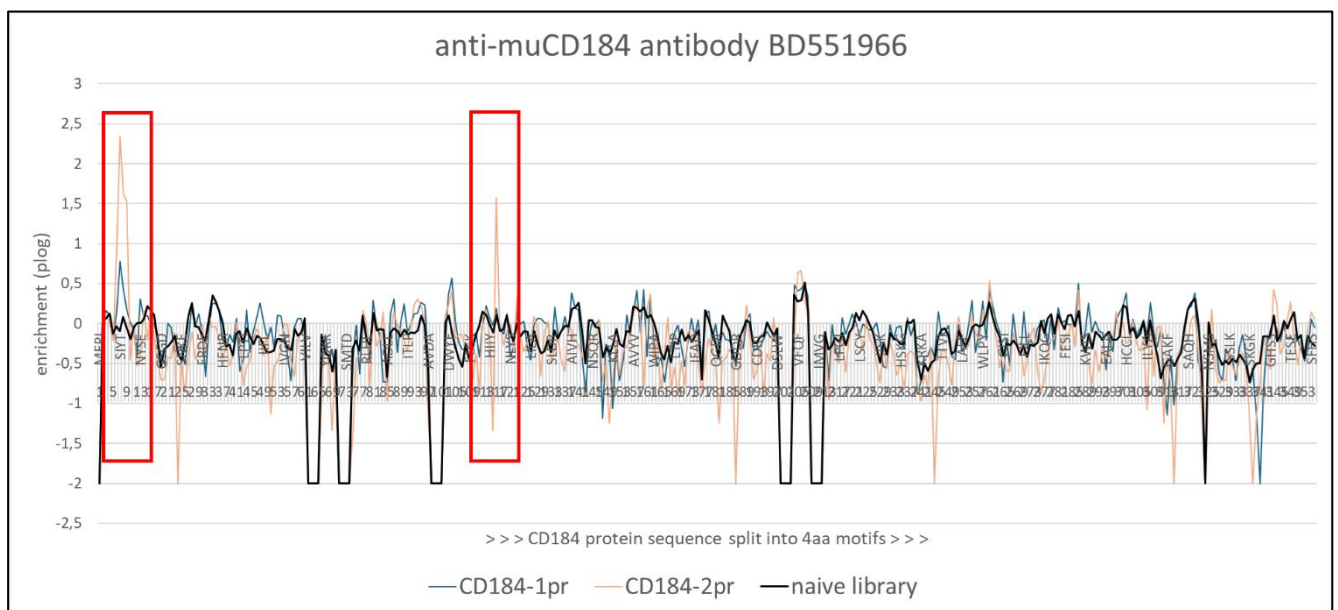
The results show/explain:

- **Epitope details by amino acids recognized**
- **Structural confirmation**
- **Species specificity**

**EPITOPIC**

## Statistical Analyses

Standard analysis based on the frequency of 4-mer motifs of the antigen sequence resulted in an overview of potential epitopes.

The graphic below shows the statistical enrichment of the consecutive 4-mer peptides of the CD184 protein sequence in the data sets after one and two selection rounds. The Y-axis shows the enrichment of 4-mers over the theoretically expected value defined by the library design. The scale is log10, so several motifs occur hundredfold and more frequently in the NGS datasets than expected. These expected values can be calculated, because the naïve library is constructed from just one codon per amino acid and behaves stable even upon replication. The black line shows the actually measured statistical data from >2 Mio. sequences of the naïve library. Deviations are due to variations in the coupling yield of codon building blocks during the synthesis.

The boxes indicate motifs within the protein enriched by more than hundredfold. Such peaks could belong to positions in the protein bound by the antibody, and the sequences containing them were further analysed.



This "enrichment" curve shows that these positions represent the motifs:
6-VSIYTSDN-13, 117-IYTV-120
It can be seen below that the second motif's frequency is due to an overall enrichment of sequences containing IYTx, where x can represent a wide variety of amino acids.

The following step in the analysis is the retrieval of full-length sequences from the selected library pool, which are containing enriched motifs. These are usually hundreds of similar sequences. They are aligned with and without the antigen sequence. In most cases the alignments were at least cured for sequences found in the datasets only once, twice or even more frequently. If a sequence is found at least three to five times, it is less likely to contain typical NGS sequencing errors.

**EPITOPIC**

## Alignment of sequences with enriched motifs

All sequences containing the amino acid 4-mers covering the partial CD184 sequence
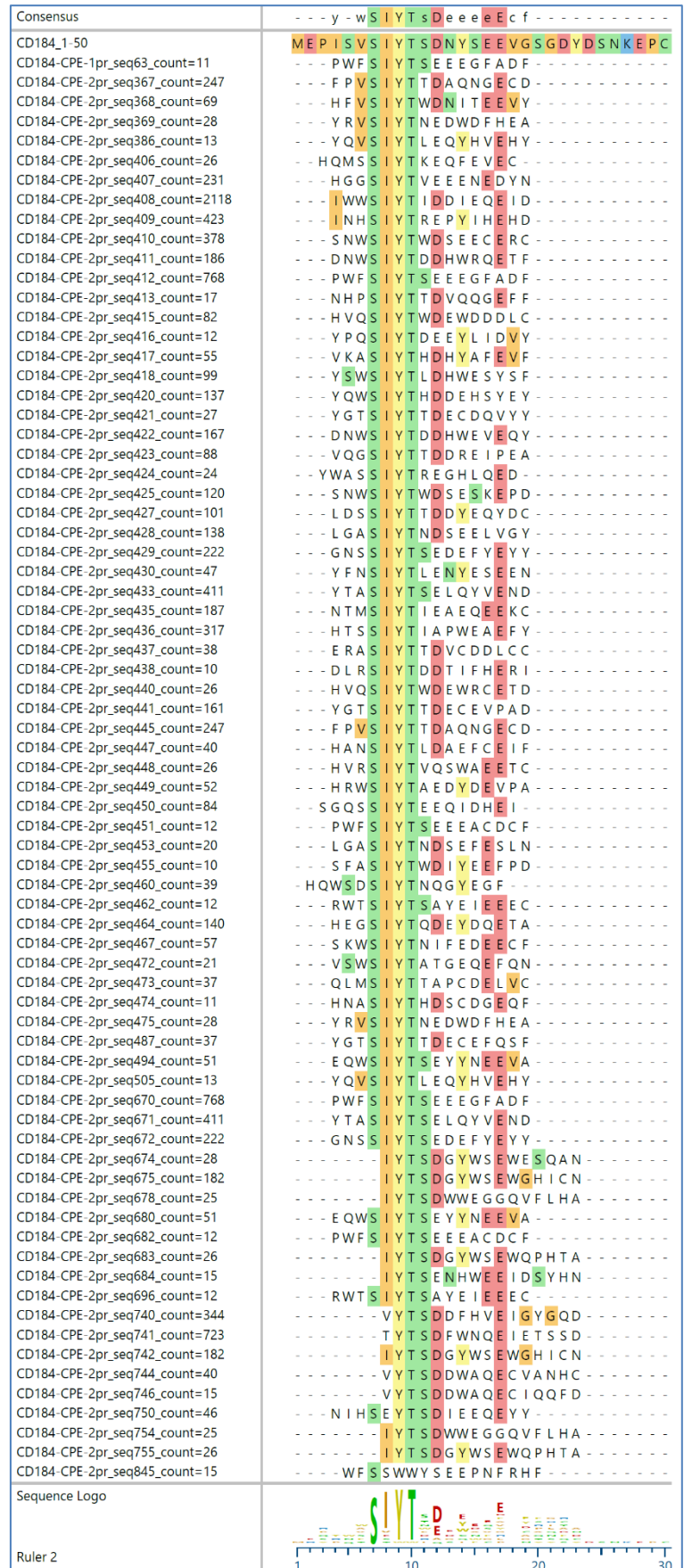
5-SVSIYTSDNYSEEVGSG-21

were retrieved from the two selection's NGS datasets. This identified 905 different in total 13,068 sequences. To keep the number of sequences displayed in the alignment to a minimum, the alignment displays only sequences found at least tenfold and sharing at least 5 amino acid identities with the aligned antigen sequence.

The epitope fingerprinting avoids multiple selection rounds. Therefore, hundreds of similar sequences can be used for the analysis of the spectrum of peptides recognized by the antibody. Due to the ability of statistical analysis of a balanced library design the selected 13,000 sequences are making up only a fraction of the 400,000 sequences of the analysed data sets.

The huge number of sequences in the ENTE-1 library (ca $5 \times 10^9$) allows also to identify multiple sequences sharing more than eight amino acid identity with the antigen. The probability to identify a **single** sequence with eight identities by chance is roughly $4 \times 10^{-11}$ (=$1/20^8$) or $2 \times 10^{-5}$ in this dataset of 500,000 sequences.

Based on the alignment and considering the amino acid residues shared by the enriched sequences the epitope could be best described as

**sIYTxDxyxeE**

| Name | Sequence |
|---|---|
| Consensus | - - - - y - w S I Y T s D e e e E c f - - - - - - - - - - |
| CD184_1-50 | M E P I S V S I Y T S D N Y S E E V G S G D Y D S N K E P C |
| CD184-CPE-1pr_seq63_count=11 | - - - P W F S I Y T S E E E G F A D F - - - - - - - - - - |
| CD184-CPE-2pr_seq367_count=247 | - - - F P V S I Y T T D A Q N G E C D - - - - |
| CD184-CPE-2pr_seq368_count=69 | - - - H F V S I Y T W D N I T E E V Y - - - |
| CD184-CPE-2pr_seq369_count=28 | - - - Y R V S I Y T N E D W D F H E A - - - |
| CD184-CPE-2pr_seq386_count=13 | - - - Y Q V S I Y T L E Q Y H V E H Y - - - |
| CD184-CPE-2pr_seq406_count=26 | - - H Q M S I Y T K E Q F E V E C - - - |
| CD184-CPE-2pr_seq407_count=231 | - - H G G S I Y T V E E E N E D Y N - - - |
| CD184-CPE-2pr_seq408_count=2118 | - - - I W W S I Y T I D D I E Q E I D - - - |
| CD184-CPE-2pr_seq409_count=423 | - - - I N H S I Y T R E P Y I H E H D - - - |
| CD184-CPE-2pr_seq410_count=378 | - - - S N W S I Y T W D S E E C E R C - - - |
| CD184-CPE-2pr_seq411_count=186 | - - - D N W S I Y T D D H W R Q E T F - - - |
| CD184-CPE-2pr_seq412_count=768 | - - - P W F S I Y T S E E E G F A D F - - - |
| CD184-CPE-2pr_seq413_count=17 | - - - N H P S I Y T T D V Q Q G E F F - - - |
| CD184-CPE-2pr_seq415_count=82 | - - - H V Q S I Y T W D E W D D D L C - - - |
| CD184-CPE-2pr_seq416_count=12 | - - - Y P Q S I Y T D E E Y L I D V Y - - - |
| CD184-CPE-2pr_seq417_count=55 | - - - V K A S I Y T H D H Y A F E V F - - - |
| CD184-CPE-2pr_seq418_count=99 | - - - Y S W S I Y T L D H W E S Y S F - - - |
| CD184-CPE-2pr_seq420_count=137 | - - - Y Q W S I Y T H D D E H S Y E Y - - - |
| CD184-CPE-2pr_seq421_count=27 | - - - Y G T S I Y T D E C D Q V Y Y - - - |
| CD184-CPE-2pr_seq422_count=167 | - - - D N W S I Y T D D H W E V E Q Y - - - |
| CD184-CPE-2pr_seq423_count=88 | - - - V Q G S I Y T D D R E I P E A - - - |
| CD184-CPE-2pr_seq424_count=24 | - - Y W A S S I Y T R E G H L Q E D - - - |
| CD184-CPE-2pr_seq425_count=120 | - - - S N W S I Y T W D S E S K E P D - - - |
| CD184-CPE-2pr_seq427_count=101 | - - - L D S S I Y T T D D Y E Q Y D C - - - |
| CD184-CPE-2pr_seq428_count=138 | - - - L G A S I Y T N D S E E L V G Y - - - |
| CD184-CPE-2pr_seq429_count=222 | - - - G N S S I Y T S E D E F Y E Y Y - - - |
| CD184-CPE-2pr_seq430_count=47 | - - - Y F N S I Y T L E N Y E S E E N - - - |
| CD184-CPE-2pr_seq433_count=411 | - - - Y T A S I Y T S E L Q Y V E N D - - - |
| CD184-CPE-2pr_seq435_count=187 | - - - N T M S I Y T I E A E Q E E K C - - - |
| CD184-CPE-2pr_seq436_count=317 | - - - H T S S I Y T I A P W E A E F Y - - - |
| CD184-CPE-2pr_seq437_count=38 | - - - E R A S I Y T T D V C D D L C C - - - |
| CD184-CPE-2pr_seq438_count=10 | - - - D L R S I Y T D D T I F H E R I - - - |
| CD184-CPE-2pr_seq440_count=26 | - - - H V Q S I Y T W D E W R C E T D - - - |
| CD184-CPE-2pr_seq441_count=161 | - - - Y G T S I Y T T D E C E V P A D - - - |
| CD184-CPE-2pr_seq445_count=247 | - - - F P V S I Y T T D A Q N G E C D - - - |
| CD184-CPE-2pr_seq447_count=40 | - - - H A N S I Y T L D A E F C E I F - - - |
| CD184-CPE-2pr_seq448_count=26 | - - - H V R S I Y T V Q S W A E E T C - - - |
| CD184-CPE-2pr_seq449_count=52 | - - - H R W S I Y T A E D Y D E V P A - - - |
| CD184-CPE-2pr_seq450_count=84 | - - S G Q S S I Y T E E Q I D H E I - - - |
| CD184-CPE-2pr_seq451_count=12 | - - - P W F S I Y T S E E E A C D C F - - - |
| CD184-CPE-2pr_seq453_count=20 | - - - L G A S I Y T N D S E F E S L N - - - |
| CD184-CPE-2pr_seq455_count=10 | - - - S F A S I Y T W D I Y E E F P D - - - |
| CD184-CPE-2pr_seq460_count=39 | - H Q W S D S I Y T N Q G Y E G F - - - |
| CD184-CPE-2pr_seq462_count=12 | - - - R W T S I Y T S A Y E I E E C - - - |
| CD184-CPE-2pr_seq464_count=140 | - - - H E G S I Y T Q D E Y D Q E T A - - - |
| CD184-CPE-2pr_seq467_count=57 | - - - S K W S I Y T N I F E D E E C F - - - |
| CD184-CPE-2pr_seq472_count=21 | - - - V S W S I Y T A T G E Q E F Q N - - - |
| CD184-CPE-2pr_seq473_count=37 | - - - Q L M S I Y T T A P C D E L V C - - - |
| CD184-CPE-2pr_seq474_count=11 | - - - H N A S I Y T H D S C D G E Q F - - - |
| CD184-CPE-2pr_seq475_count=28 | - - - Y R V S I Y T N E D W D F H E A - - - |
| CD184-CPE-2pr_seq487_count=37 | - - - Y G T S I Y T D E C E F Q S F - - - |
| CD184-CPE-2pr_seq494_count=51 | - - - E Q W S I Y T S E Y Y N E E V A - - - |
| CD184-CPE-2pr_seq505_count=13 | - - - Y Q V S I Y T L E Q Y H V E H Y - - - |
| CD184-CPE-2pr_seq670_count=768 | - - - P W F S I Y T S E E E G F A D F - - - |
| CD184-CPE-2pr_seq671_count=411 | - - - Y T A S I Y T S E L Q Y V E N D - - - |
| CD184-CPE-2pr_seq672_count=222 | - - - G N S S I Y T S E D E F Y E Y Y - - - |
| CD184-CPE-2pr_seq674_count=28 | - - - - - - I Y T S D G Y W S E W E S Q A N - - |
| CD184-CPE-2pr_seq675_count=182 | - - - - - - - I Y T S D G Y W S E W G H I C N - - |
| CD184-CPE-2pr_seq678_count=25 | - - - - - - - I Y T S D W W E G G Q V F L H A - - |
| CD184-CPE-2pr_seq680_count=51 | - - - E Q W S I Y T S E Y Y N E E V A - - - |
| CD184-CPE-2pr_seq682_count=12 | - - - P W F S I Y T S E E E A C D C F - - - |
| CD184-CPE-2pr_seq683_count=26 | - - - - - - I Y T S D G Y W S E W Q P H T A - - |
| CD184-CPE-2pr_seq684_count=15 | - - - - - - - I Y T S E N H W E E I D S Y H N - - |
| CD184-CPE-2pr_seq696_count=12 | - - - - - R W T S I Y T S A Y E I E E C - - - |
| CD184-CPE-2pr_seq740_count=344 | - - - - - - V Y T S D D F H V E I G Y G Q D - - |
| CD184-CPE-2pr_seq741_count=723 | - - - - - - - T Y T S D F W N Q E I E T S S D - - |
| CD184-CPE-2pr_seq742_count=182 | - - - - - - I Y T S D G Y W S E W G H I C N - - |
| CD184-CPE-2pr_seq744_count=40 | - - - - - - V Y T S D D W A Q E C V A N H C - - |
| CD184-CPE-2pr_seq746_count=15 | - - - - - - V Y T S D D W A Q E C I Q Q F D - - |
| CD184-CPE-2pr_seq750_count=46 | - - - N I H S E Y T S D I E E Q E Y Y - - - |
| CD184-CPE-2pr_seq754_count=25 | - - - - - - I Y T S D W W E G G Q V F L H A - - |
| CD184-CPE-2pr_seq755_count=26 | - - - - - - I Y T S D G Y W S E W Q P H T A - - |
| CD184-CPE-2pr_seq845_count=15 | - - - - W F S S W W Y S E E P N F R H F - - - |
| Sequence Logo | (sequence logo graphic) |
| Ruler 2 | 1 ........ 10 ........ 20 ........ 30 |

The amino acid sequence of the epitope is not a continuous stretch. Therefore, it makes sense to search for discontinuous motifs, too. The following result is obtained when sequences are retrieved from the data base containing 5-mer variations of the type AxAAA, AAxAA, AAAxA of the motifs SVSIY, SIYTS, YTSDN, SDNYS or NYSEE (taken from 5-SVSIYTSDNYSEEVGSG-21). Here x can be any amino acid.

This approach can possibly identify more sequences by allowing variations of amino acids instead of strict 4-mer motifs. In this case, it expands the alignment above by only showing the sequences which are found at least 5 times and that have 5 or more amino acids identity to the antigen sequence. But the sequences in total only repeat the pattern which was visible in the previous alignment already.

Based on the second alignment the epitope would be best described as below. Here for example xYx can be replaced by one or two other hydrophobic amino acids and a high probability for two Glu being recognized by the antibody.

# vsIYTSdxYxEEv

In general such amino acids must not always be interacting with the antibody, but they might also be essential for a proper fold of the peptide. One example could be the Asp in the center of the sequence.

```
Ruler 1                                  1        10        20
Consensus                                - - - y w s S I Y T S d y y e E e d - - -

CD184_1-55                               M E P I S V S I Y T S D N Y S E E V G S G D
seq56_CD184-CPE-2pr_count_69             - - - H F V S I Y T W D N I T E E V Y - - -
seq69_CD184-CPE-2pr_count_247            - - - - F P V S I Y T T D A Q N G E C D - - -
seq71_CD184-CPE-2pr_count_13             - - - Y Q V S I Y T L E Q Y H V E H Y - - -
seq73_CD184-CPE-2pr_count_7              - - - P E V S I Y T Q D R I G P E Q D - - -
seq76_CD184-CPE-2pr_count_9              - - - V N V S I Y E I D D W C E E R D - - -
seq111_CD184-CPE-2pr_count_21            - - - V S W S I Y T A T G E Q E F Q N - - -
seq118_CD184-CPE-2pr_count_8             - R Y V S Y S I Y T N E H I - - F E D - - -
seq135_CD184-CPE-2pr_count_39            - H Q W S D S I Y T N Q G Y E G F - - - - -
seq137_CD184-CPE-2pr_count_99            - - - Y S W S I Y T L D H W E S Y S F - - -
seq232_CD184-CPE-2pr_count_89            - - - - S S V S V Y T N E E Q H A E V C - - -
seq233_CD184-CPE-2pr_count_7             - - - - N S V S D Y T G D L Q E E E L Y - - -
seq245_CD184-CPE-2pr_count_8             - - - W S V S F Y S K D A P D D E C A - - -
seq252_CD184-CPE-2pr_count_6             - - - V S V S W S I Y T E Q Y V E N D - - -
seq2_CD184-CPE-1pr_count_5               - - - Y T A S I Y T S E L Q Y V E N D - - -
seq7_CD184-CPE-2pr_count_7               - E Q W S - I Y T S E Y Y D L T E A - - -
seq8_CD184-CPE-2pr_count_411             - - Y T A S I Y T S E L Q Y V E N D - - -
seq9_CD184-CPE-2pr_count_222             - - - G N S S I Y T S E D E F Y E Y Y - - -
seq27_CD184-CPE-2pr_count_51             - - E Q W S - I Y T S E Y Y N E E V A - - -
seq37_CD184-CPE-2pr_count_12             - - - R W T S I Y T S A Y E I E E E C - - -
seq39_CD184-CPE-1pr_count_5              - - - Y T A S I Y T S E L Q Y V E N D - - -
seq53_CD184-CPE-1pr_count_5              - - - - - - - I Y T S D G Y W S E W G H I C
seq79_CD184-CPE-2pr_count_411            - - - Y T A S I Y T S E L Q Y V E N D - - -
seq100_CD184-CPE-2pr_count_51            - - E Q W S - I Y T S E Y Y N E E V A - - -
seq106_CD184-CPE-2pr_count_182           - - - - - - - I Y T S D G Y W S E W G H I C
seq112_CD184-CPE-2pr_count_222           - - - G N S S I Y T S E D E F Y E Y Y - - -
seq122_CD184-CPE-2pr_count_26            - - - - - - - I Y T S D G Y W S E W Q P H T
seq128_CD184-CPE-2pr_count_15            - - - - - - - I Y T S E N H W E E I D S Y H
seq135_CD184-CPE-2pr_count_12            - - - R W T S I Y T S A Y E I E E E C - - -
seq139_CD184-CPE-2pr_count_28            - - - - - - - I Y T S D G Y W S E W E S Q A
seq140_CD184-CPE-2pr_count_7             - - E Q W S - I Y T S E Y Y D L T E A - - -
seq150_CD184-CPE-1pr_count_5             - - - Y T A S I Y T S E L Q Y V E N D - - -
seq182_CD184-CPE-2pr_count_9             - - - R Q M S H Y T S D D P V L E D C - - -
seq184_CD184-CPE-2pr_count_222           - - - G N S S I Y T S E D E F Y E Y Y - - -
seq185_CD184-CPE-2pr_count_51            - - E Q W S - I Y T S E Y Y N E E V A - - -
seq189_CD184-CPE-2pr_count_12            - - - R W T S I Y T S A Y E I E E E C - - -
seq220_CD184-CPE-2pr_count_20            - - - R P V S V Y T S V Q E V D E G F - - -
seq226_CD184-CPE-2pr_count_23            - - - T A N S Q Y T S S N Y H D E T Y - - -
seq227_CD184-CPE-2pr_count_46            - - N I H S E Y T S D I E E Q E Y Y - - -
seq230_CD184-CPE-2pr_count_6             - - - - - I L S S Y T S D F E C E E F V A -
seq237_CD184-CPE-2pr_count_19            - - - I H A S N Y T S A P Y Y K E Q D - - -
seq243_CD184-CPE-2pr_count_14            - - - I W Q S D Y T S E D Y Q S E R Y - - -
seq248_CD184-CPE-2pr_count_6             - - - - A V V S E Y T S E P Y D Q E K D - - -
seq253_CD184-CPE-2pr_count_33            - - - L W N S T Y T S H S Q Q E E H D - - -
seq256_CD184-CPE-2pr_count_7             - E Q W S - I Y T S E Y Y D L T E A - - -
seq260_CD184-CPE-2pr_count_6             - - - - - E H S S Y T S D A F E Q D Q G C - -
seq273_CD184-CPE-2pr_count_411           - - Y T A S I Y T S E L Q Y V E N D - - -
seq281_CD184-CPE-2pr_count_7             - - Q T M S E Y T S D T C Q Q E V F - - -
seq282_CD184-CPE-2pr_count_5             - - - G E H S Q Y T S Q T Y F E E T A - - -
seq301_CD184-CPE-1pr_count_5             - - - Y T A S I Y T S E L Q Y V E N D - - -
seq304_CD184-CPE-2pr_count_411           - - - Y T A S I Y T S E L Q Y V E N D - - -
seq316_CD184-CPE-2pr_count_51            - - E Q W S - I Y T S E Y Y N E E V A - - -
seq325_CD184-CPE-2pr_count_222           - - - G N S S I Y T S E D E F Y E Y Y - - -
seq337_CD184-CPE-2pr_count_12            - - - R W T S I Y T S A Y E I E E E C - - -
seq338_CD184-CPE-2pr_count_5             - - - E I M S I W T S I E Y L S E A D - - -
seq342_CD184-CPE-2pr_count_7             - E Q W S - I Y T S E Y Y D L T E A - - -
seq431_CD184-CPE-1pr_count_5             - - - Y T A S I Y T S E L Q Y V E N D - - -
seq438_CD184-CPE-2pr_count_222           - - - G N S S I Y T S E D E F Y E Y Y - - -
seq439_CD184-CPE-2pr_count_51            - - E Q W S - I Y T S E Y Y N E E V A - - -
seq442_CD184-CPE-2pr_count_12            - - - R W T S I Y T S A Y E I E E E C - - -
seq515_CD184-CPE-2pr_count_7             - - E Q W S - I Y T S E Y Y D L T E A - - -
seq539_CD184-CPE-2pr_count_411           - - - Y T A S I Y T S E L Q Y V E N D - - -
seq621_CD184-CPE-1pr_count_5             - - - Y T A S I Y T S E L Q Y V E N D - - -
seq633_CD184-CPE-2pr_count_51            - - E Q W S - I Y T S E Y Y N E E V A - - -
seq644_CD184-CPE-2pr_count_7             - - - P E V S I Y T Q D R I G P E Q D - - -
seq646_CD184-CPE-2pr_count_7             - - - V Y R S I Y T L D Y R Q L E L N - - -
seq649_CD184-CPE-2pr_count_69            - - - H F V S I Y T W D N I T E E V Y - - -
seq650_CD184-CPE-2pr_count_411           - - - Y T A S I Y T S E L Q Y V E N D - - -
seq654_CD184-CPE-2pr_count_6             - - - V L S S I Y T V D Q Y E Q G L I - - -
seq673_CD184-CPE-2pr_count_187           - - N T M S I Y T I E A E Q E E K C - - -
seq681_CD184-CPE-2pr_count_37            - - - Q L M S I Y T T A P C D E - L V C - -
seq685_CD184-CPE-2pr_count_423           - - - I N H S I Y T R E P Y I H E H D - - -
seq686_CD184-CPE-2pr_count_12            - - - Y P Q S I Y T D E E Y L I D V Y - - -
seq692_CD184-CPE-2pr_count_6             - - - W P P S I Y T R D R E D E K L N - - -
seq700_CD184-CPE-2pr_count_52            - - - H R W S I Y T A E D Y D E V P A - - -
seq704_CD184-CPE-2pr_count_140           - - - H E G S I Y T Q D E Y D Q E T A - - -
seq713_CD184-CPE-2pr_count_39            - H Q W S D S I Y T N Q G Y E G F - - - - -
seq724_CD184-CPE-2pr_count_12            - - - R W T S I Y T S A Y E I E E E C - - -
seq729_CD184-CPE-2pr_count_247           - - - - F P V S I Y T T D A Q N G E C D - - -
seq730_CD184-CPE-2pr_count_8             - R Y V S Y S I Y T N E H I - - F E D - - -
seq737_CD184-CPE-2pr_count_10            - - - D L R S I Y T D D T I F H E R I - - -
seq750_CD184-CPE-2pr_count_222           - - - G N S S I Y T S E D E F Y E Y Y - - -
seq758_CD184-CPE-2pr_count_2118          - - - I W W S I Y T I D D I E Q E I D - - -
seq765_CD184-CPE-2pr_count_11            - - - H N A S I Y T H D S C D G E Q F - - -
seq766_CD184-CPE-2pr_count_101           - - - L D S S I Y T T D D Y E Q Y D C - - -
seq767_CD184-CPE-2pr_count_167           - - - D N W S I Y T D D H W E V E Q Y - - -
seq769_CD184-CPE-2pr_count_186           - - - D N W S I Y T D D H W R Q E T F - - -
seq775_CD184-CPE-2pr_count_7             - - E Q W S - I Y T S E Y Y D L T E A - - -
seq776_CD184-CPE-2pr_count_26            - - - H V R S I Y T V Q S W A E E T C - - -
seq800_CD184-CPE-2pr_count_21            - - - V S W S I Y T A T G E Q E F Q N - - -
seq806_CD184-CPE-2pr_count_13            - - - Y Q V S I Y T L E Q Y H V E H Y - - -
seq809_CD184-CPE-2pr_count_5             - - - A V W S I Y T L E D Q E V I T I - - -
seq810_CD184-CPE-2pr_count_99            - - Y S W S I Y T L D H W E S Y S F - - -
seq814_CD184-CPE-2pr_count_26            - - - H V Q S I Y T W D E W R C E T D - - -
seq815_CD184-CPE-2pr_count_57            - - - S K W S I Y T N I F E D E E C F - - -
seq829_CD184-CPE-2pr_count_47            - - - Y F N S I Y T L E N Y E S E E N - - -
seq849_CD184-CPE-2pr_count_10            - - - S F A S I Y T W D I - - Y E E F P D -
seq851_CD184-CPE-2pr_count_120           - - - S N W S I Y T W D S - - E S K E P D -
seq854_CD184-CPE-2pr_count_20            - - - L G A S I Y T N D S E F E S L N - - -
seq858_CD184-CPE-2pr_count_5             - - - Y K W S I Y T W D P I H C E T I - - -
seq862_CD184-CPE-2pr_count_40            - - - H A N S I Y T L D A E F C E I F - - -
seq876_CD184-CPE-2pr_count_17            - - - N H P S I Y T T D V Q Q G E F F - - -
seq877_CD184-CPE-2pr_count_55            - - - V K A S I Y T H D H Y A F E V F - - -
seq884_CD184-CPE-2pr_count_378           - - - S N W S I Y T W D S - - E E C E R C -
seq71_CD184-CPE-2pr_count_26             - - - - I H S S Y T T D N C E E S R F C -
seq76_CD184-CPE-2pr_count_69             - - - H F V S I Y T W D N I T E E V Y - - -
seq80_CD184-CPE-2pr_count_28             - - - R H H S E Y T V D N E H Q E R A - - -
seq81_CD184-CPE-2pr_count_6              - - - N I P S N Y T E D N Y A D E F N - - -
seq113_CD184-CPE-2pr_count_23            - - - T A N S Q Y T S S N Y H D E T Y - - -
seq120_CD184-CPE-2pr_count_15            - - - - - - - I Y T S E N H W E E I D S Y H
seq133_CD184-CPE-1pr_count_5             - - - - - - - I Y T S D G Y W S E W G H I C
seq148_CD184-CPE-2pr_count_9             - - - R Q M S H Y T S D D P V L E D C - - -
seq163_CD184-CPE-2pr_count_182           - - - - - - - I Y T S D G Y W S E W G H I C
seq165_CD184-CPE-2pr_count_344           - - - - - - - V Y T S D D F H V E I G Y G Q
seq174_CD184-CPE-2pr_count_46            - - N I H S E Y T S D I E E Q E Y Y - - -
seq176_CD184-CPE-2pr_count_6             - - - - - I L S S Y T S D F E C E E F V A -
seq194_CD184-CPE-2pr_count_7             - - - - - - - T Y T S D H E D E E Y H G Y N
seq198_CD184-CPE-2pr_count_6             - - - E H S S Y T S D A F E Q D Q G C - -
seq202_CD184-CPE-2pr_count_9             - - - - - - - I Y T S D G Y W S E W Q P H T
```

**EPITOPIC**

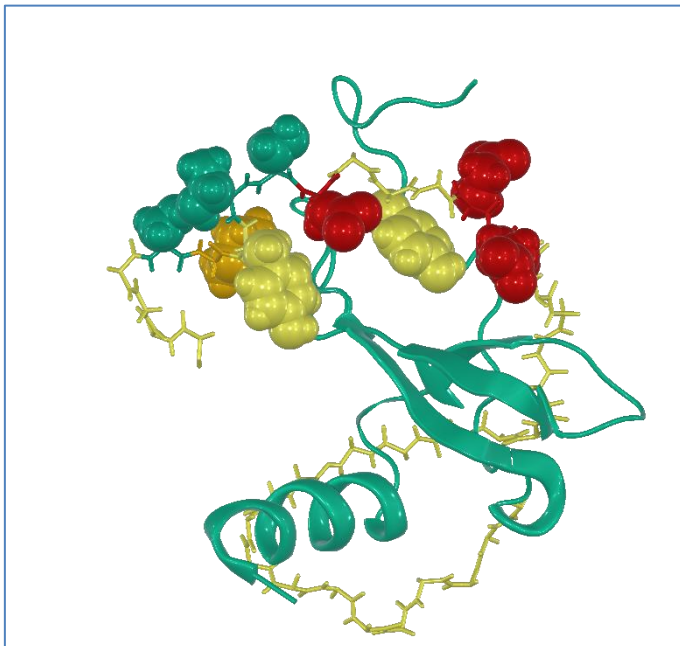## CD184 epitope location in the protein sequence

This is an overview of the full-length CD184 sequence where those regions are marked in yellow, that show enrichment in the statistical analysis of 4-mer motifs in the NGS datasets.

```
> CXCR4_MOUSE C-X-C chemokine receptor type 4 (P70658)

1       MEPISVSIYT SDNYSEEVGS GDYDSNKEPC FRDENVHFNR IFLPTIYFII FLTGIVGNGL
61      VILVMGYQKK LRSMTDKYRL HLSVADLLFV ITLPFWAVDA MADWYFGKFL CKAVHIIYTV
121     NLYSSVLILA FISLDRYLAI VHATNSQRPR KLLAEKAVYV GVWIPALLLT IPDFIFADVS
181     QGDISQGDDR YICDRLYPDS LWMVVFQFQH IMVGLILPGI VILSCYCIII SKLSHSKGHQ
241     KRKALKTTVI LILAFFACWL PYYVGISIDS FILLGVIKQG CDFESIVHKW ISITEALAFF
301     HCCLNPILYA FLGAKFKSSA QHALNSMSRG SSLKILSKGK RGGHSSVSTE SESSSFHSS
```

## CD184 epitope in structures

CD184 or CXCR4 is a transmembrane receptor and the antibody recognizes the extracellular N-terminus. Therefore, no structures of the full protein have been published. The picture below shows the otherwise unstructured N-terminal sequence (yellow) in complex with a ligand. **SIYTSDnYsEE** side chains are displayed as spheres, the ligand as green cartoon. This structure (PDB 2k04) shows the N-terminus of the human protein. It can be seen how the amino acids forming the epitope could be accessed by  an antibody. This structural view explains also, why the epitope is not a continuous sequence. The specificity of this antibody results not from a defined structure, but from the number of conserved residues recognized in a (without ligand) disordered structure.

**EPITOPIC**

### Species specificity

The antibody originating from rat is well described as specific for the mouse protein. Nevertheless, the N-terminus of the CXCR4 is well-conserved among different species.



According to this alignment it is most likely, that N- and C-terminus of the epitope recognized by the antibody are important for this specificity, while the central region is shared amongst the different receptor proteins.

**Contact:**

**epitopic GmbH – Deutscher Platz 5e – 04103 Leipzig – Germany**

**E-Mail: info@epitopic.com**

**Phone +49 341 253 55 160**

**www.epitopic.com**