**Identification of candidate transcytosis proteins through sequence analysis.**
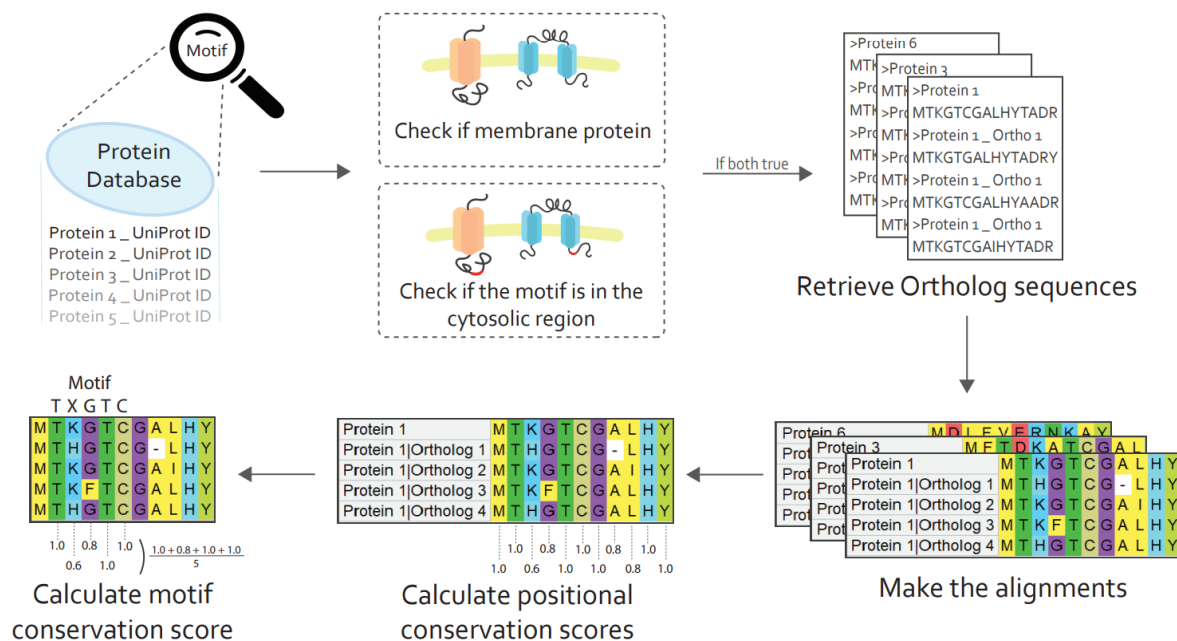
A set of motifs has been identified from the literature, which has been demonstrated to play a crucial role in the transcytosis process (as detailed in Table 1). To further investigate these motifs, the REST (Representational State Transfer) web service of the PrositeScan tool (de Castro et al., 2006) provided by the ExPASY portal was employed to search for proteome sequences that possess these motifs. The UniProt proteome database was searched for each motif, with isoforms filtered out and the search restricted to the Homo Sapiens organism. Metadata is generated using UniProt REST API. The metadata includes the location of the motifs on the sequence, topology information, and database reference IDs of each protein. Through the analysis of topology information, proteins without transmembrane domains are filtered out. The proteins that contain the motif in their cytoplasmic domains were selected from the remaining protein sequences. These proteins have been designated as preliminary candidates. The orthologous protein sequences of the preliminary candidate proteins are obtained from the OMA (Orthologous Matrix) (Altenhoff vd., 2021) database. Using the orthologous sequences, multiple sequence alignments (MSAs) were built for each pre-candidate protein using the Clustal Omega (Sievers vd., 2011) alignment method.

A custom Python script was used to calculate the positional conservation scores, which are determined by the ratio of the most common amino acid at that position to the total number of sequences in the alignment. The average conservation score of each MSA has been calculated by taking the mean conservation score of the positions for the alignment. While taking this average, the positions where the most observed character is a gap have been noted as significant contributors to the mean of the MSA, specifically for the sequences with more distant common ancestors. Hence, two mean MSA conservation scores have been generated, both with and without the positions primarily composed of gaps. The mean motif score has been calculated by averaging the conservation scores of the motif positions. For the motifs which have an X character, the X characters conservation score has been neglected for the calculation of the mean motif score. The percentile of the motif's conservation score is determined from the data by a custom Python script in order to assess the degree of preservation of the motif region in comparison to the rest of the MSA. The percentile of the motifall positions
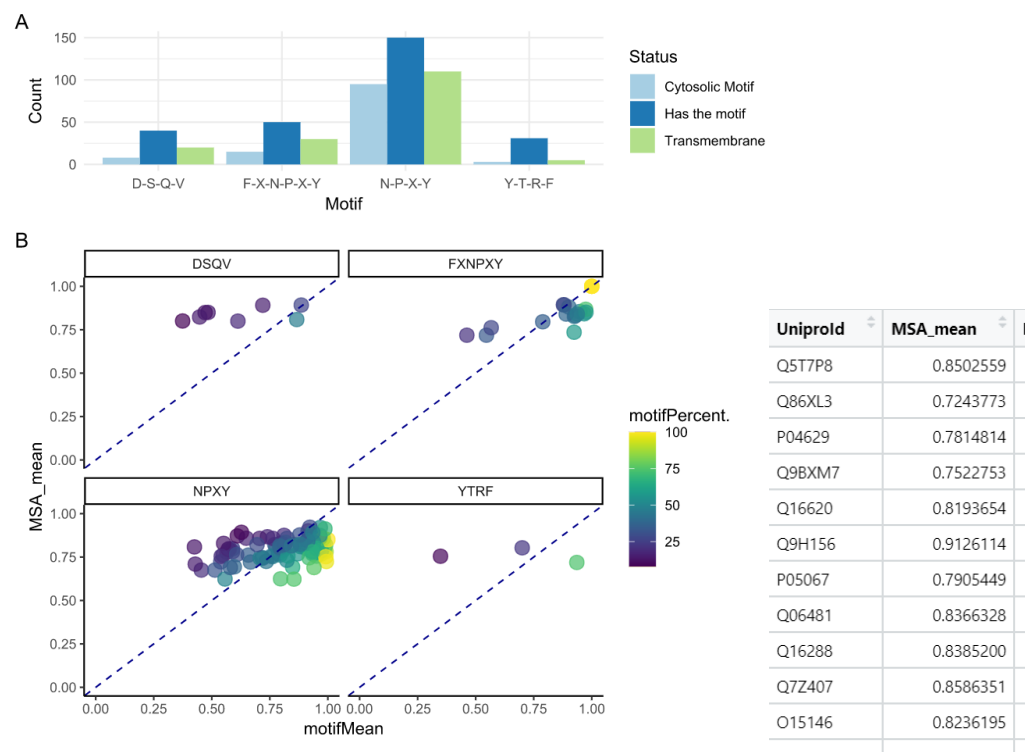


Figure X:

After the filtration of non-transmembrane proteins and motif location, the pre-candidate counts for each motif have been found as; 3 for YTRF, 109 for NPXY, 19 for FXNPXY, and 8 for DSQV. 82 of the NPXY protein showed a higher motif mean than the mean of their MSA. Out of these proteins, 10 of them have >98% conservation score.
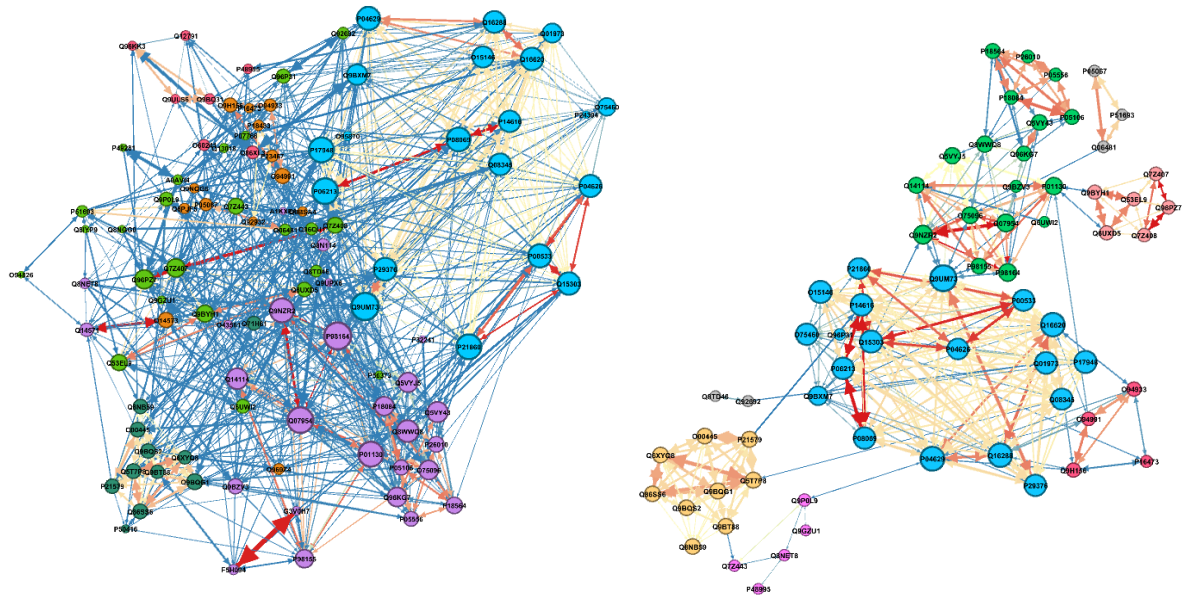
Figure X :