

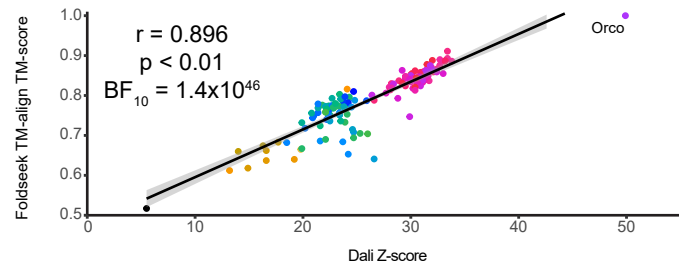
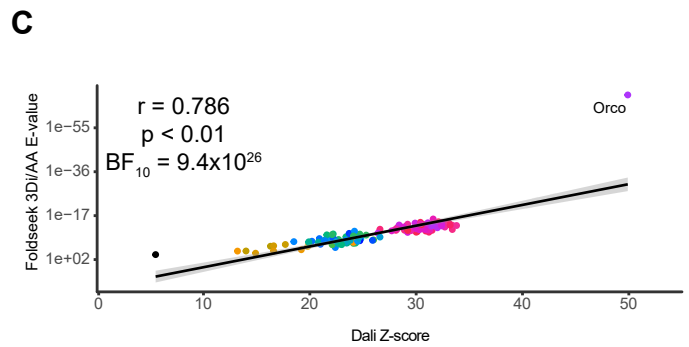
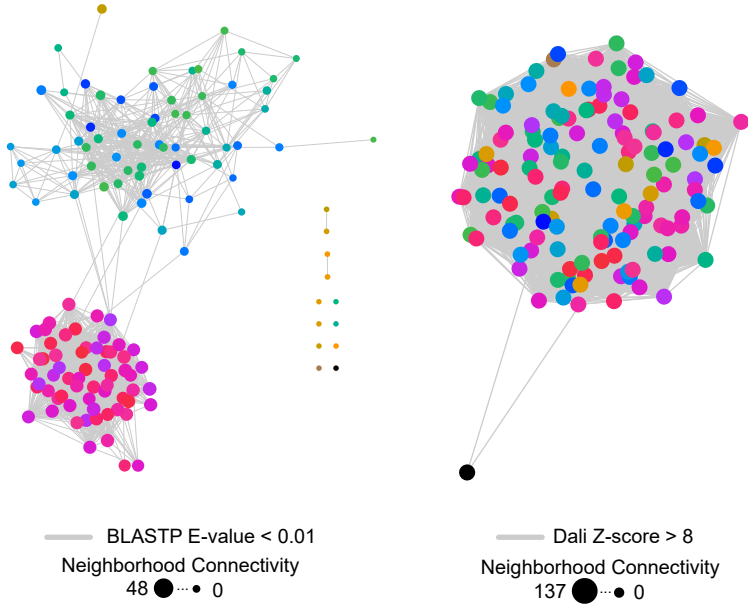
**Current Biology, Volume 33**

**Supplemental Information**

**Remote homolog detection places insect  
chemoreceptors in a cryptic protein superfamily  
spanning the tree of life**

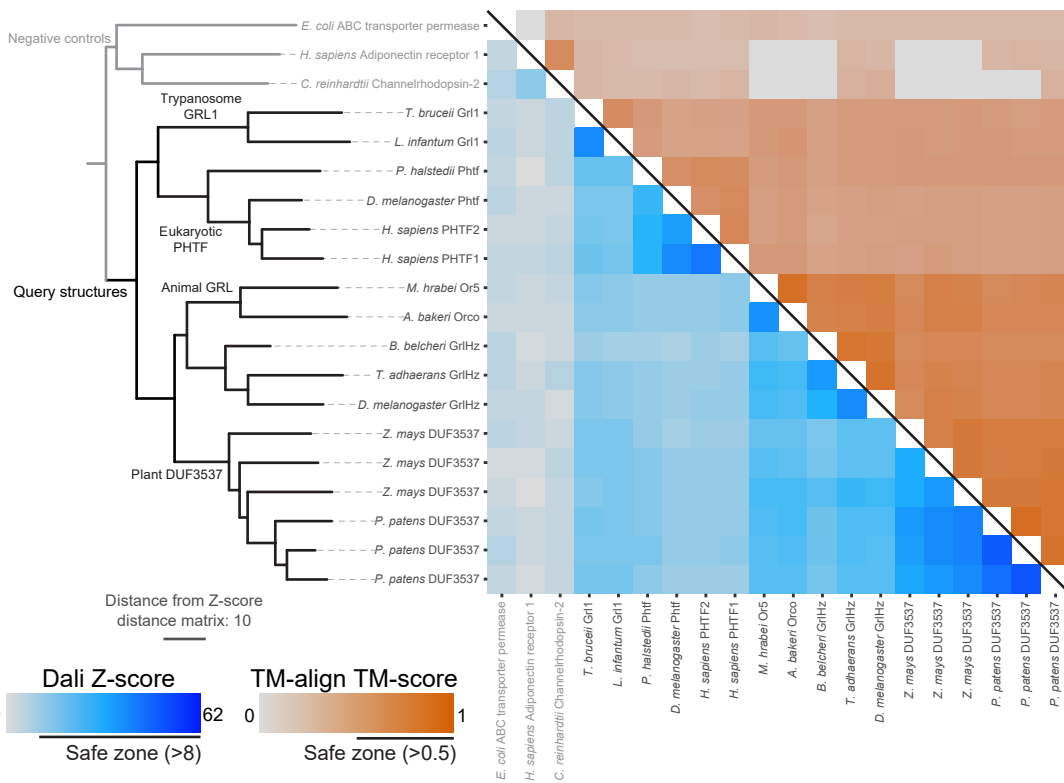
**Nathaniel J. Himmel, David Moi, and Richard Benton**

**A** Sequence similarity network **B** Structure similarity network

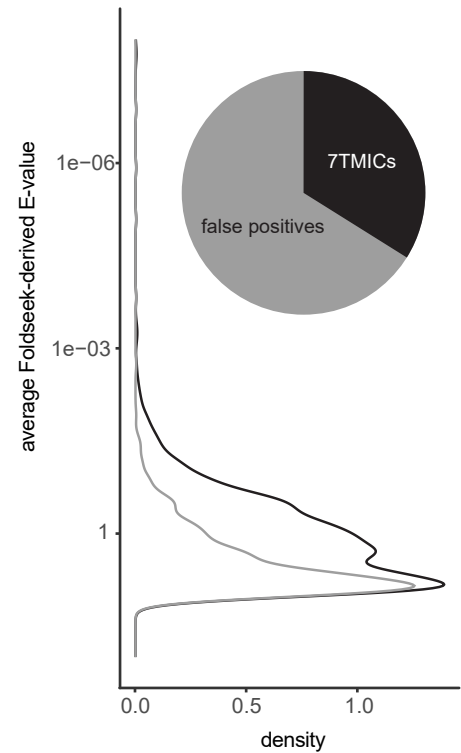


● Phtf ● GrHz ● Grs ● Ors ● Grls

**D**



**E**



**Figure S1. All-to-all pairwise protein similarity networks of *D. melanogaster* 7TMICs, Foldseek benchmarking, and summary of the Foldseek screen, related to Figure 1 and Figure 2.**

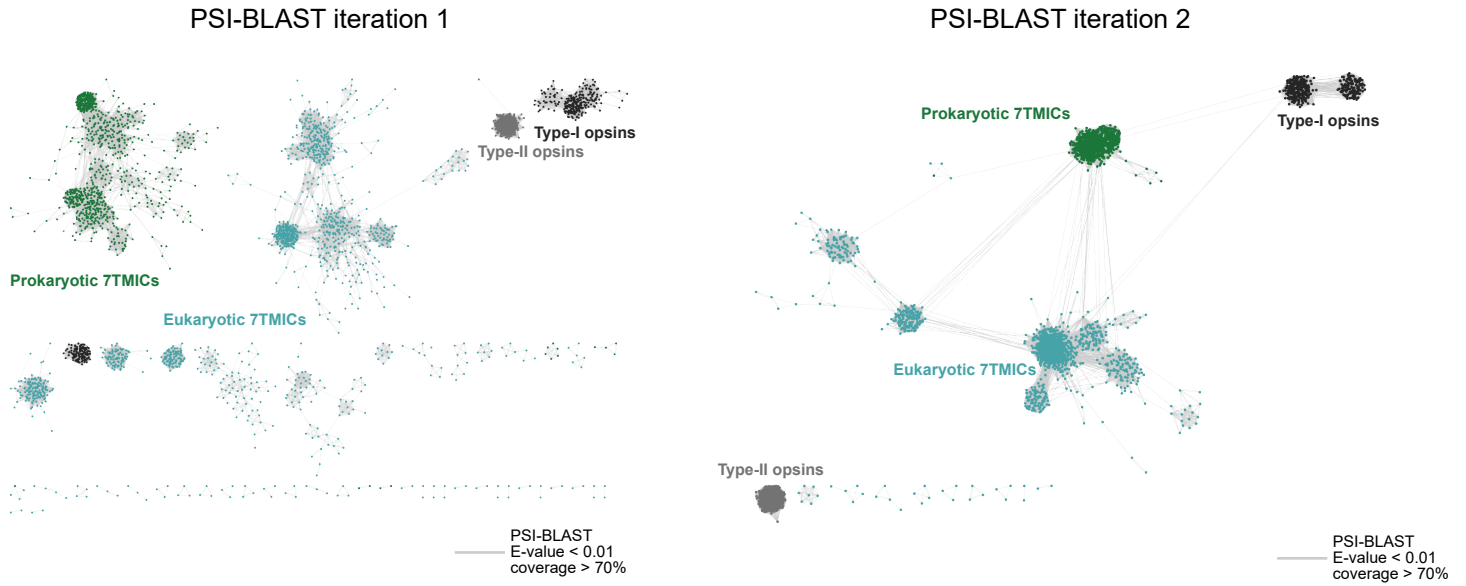
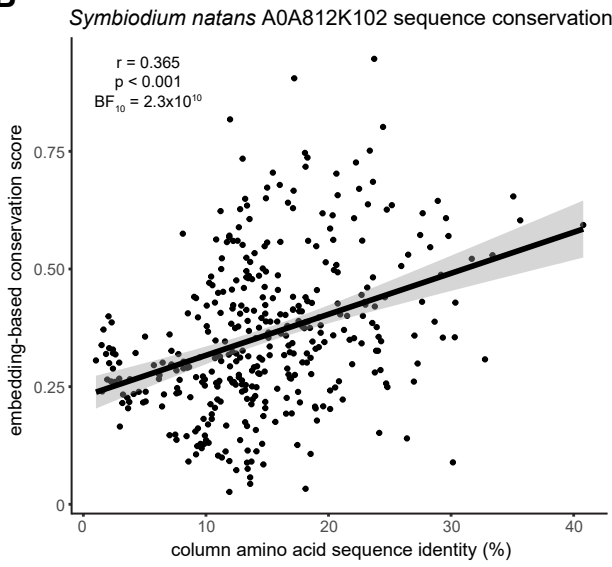
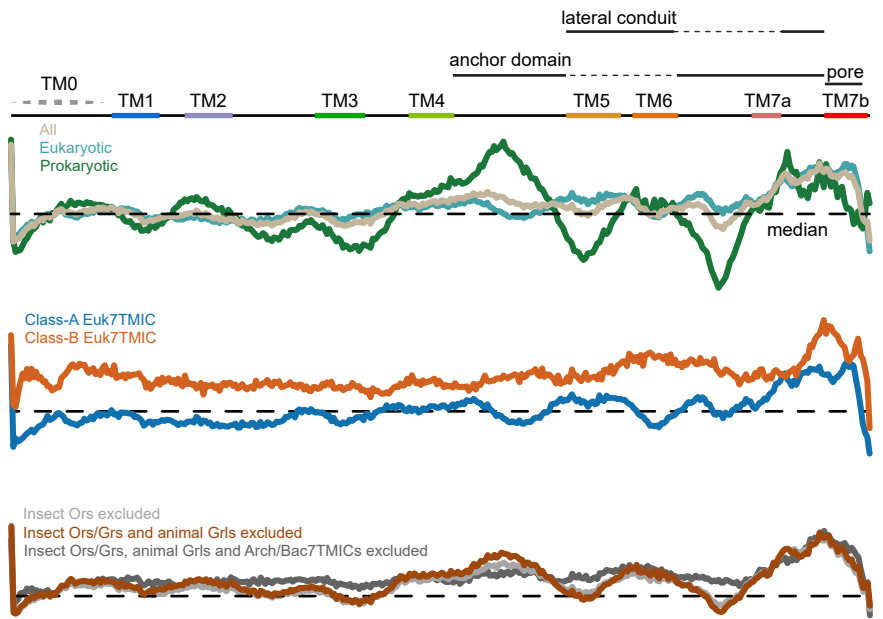
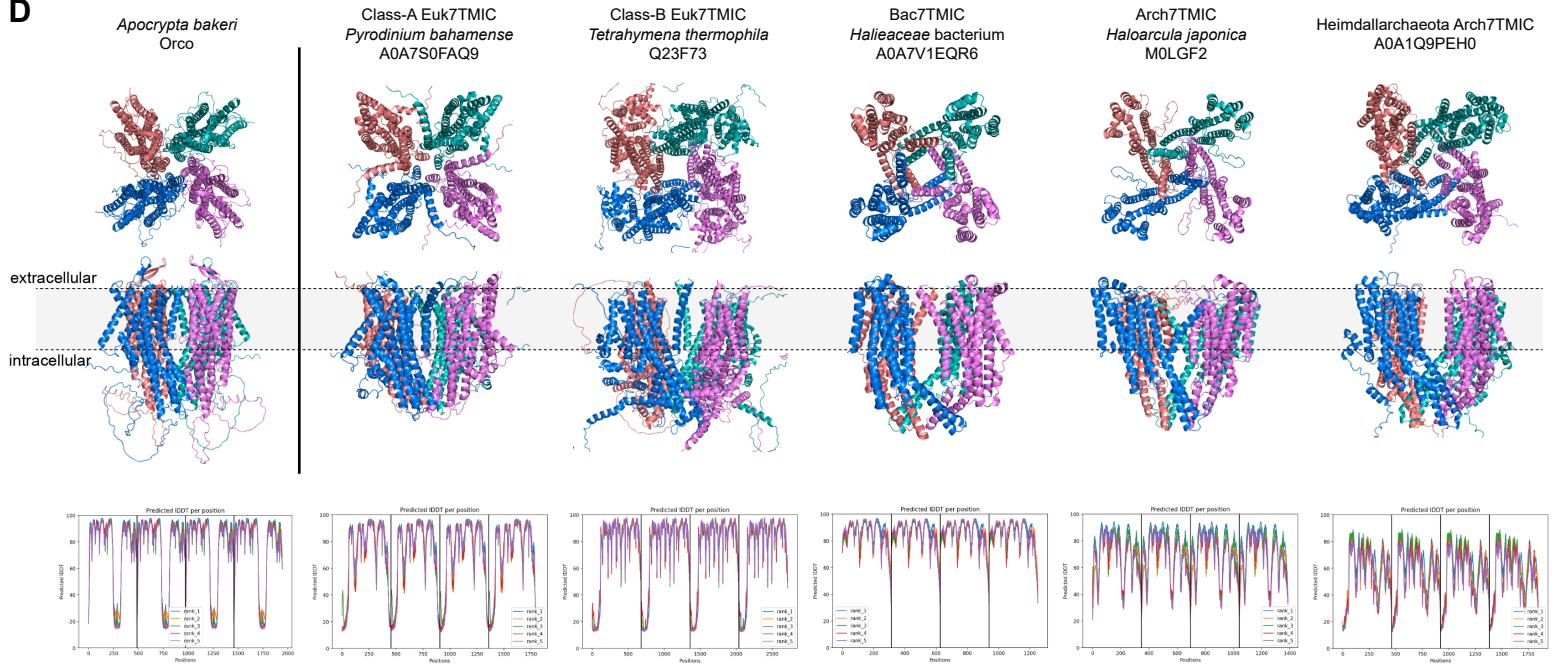
(A) All-to-all BLASTP network of *D. melanogaster* 7TMICs; consistent with their low pairwise sequence similarity, this analysis fails to link every 7TMIC to all others. Rather, the major *D. melanogaster* classes (Ors and Grs) are separated into two identifiable community structures, with sparse connectivity among the Grs, and between the Grs and Ors. Other 7TMICs—including GrIs, GrHz, Phtf and two Grs—form singlets, indicating an inability to identify hypothetical homologs using BLASTP. The color key is shown below panels (A-C), and matches that of **Figure 1E**.

(B) All-to-all Dali network of *D. melanogaster* 7TMICs. In contrast to (A), structural comparisons result in a “hairball” network, where nearly all proteins are linked to all others, excepting Phtf, which is presumed to be the most distantly related.

(C) Plots of structural similarity scores between Orco and other *D. melanogaster* 7TMICs, comparing Dali to Foldseek-derived scores. Foldseek generates Orco-to-all E-values that tightly correlate with the rapidly generated 3Di+AA-derived E-values (top) and the slowly generated TM-align-derived TM-scores (bottom).

(D) Protein models used in the Foldseek screen, and negative controls used for subsequent Dali-based validation, with a clustering dendrogram based on all-to-all Dali comparisons between the queries and negative controls. The dendrogram is derived from the Dali Z-score distance matrix. The annotation on the heatmap corresponds to the “groups” described in the methods. The heatmap shows all-to-all Dali Z-scores and TM-scores. All 7TMIC-to-control comparisons are well below thresholds of confidence in fold similarity: the Dali Z-score maximum is 7.1 and averages 4.3; and the TM-align TM-score maximum is 0.31 and averages 0.2. By contrast, for 7TMIC-to-7TMIC comparisons: the Z-score minimum is 8.5 and averages 22.0, and the TM-score minimum is 0.4 and averages 0.6.

(E) Stacked density plot showing the frequency distribution of the hits of the Foldseek screen, by E-value, with the inset pie-chart showing the proportion of true positives to false positives. Many true positives had relatively poor E-values, with similar or worse scores than many false positives, demonstrating the need for structural validation in a Foldseek screen.

**A****B****C****D**

**Figure S2. Initial iterations of the PSI-BLAST sequence similarity networks, 7TMIC sequence conservation analysis, and predicted quaternary structures of select, newly-identified 7TMICs, related to Figure 3.**

(A) Sequence similarity networks were generated by all-to-all PSI-BLAST searches of a 50% clustered sequence database of 7TMICs, alongside databases of Type-I and Type-II opsins. Iterations 1 and 2 are visualized here. Subsequent iterations resemble the clustering pattern of iteration 3, as visualized in **Figure 3B**, albeit with strengthening community structures. Left: PSI-BLAST iteration 1. In this network, sequences formed several non-contiguous clusters, and failed to cluster together 7TMICs and Type-I opsins, which is expected given the substantial sequence dissimilarity of 7TMICs. Right: PSI-BLAST iteration 2. Surprisingly, PSI-BLAST networking produced bidirectional linking of the majority of 7TMICs, although presumed spurious linkages to outgroups began to form (which did not greatly multiply in subsequent iterations), and a small number of 7TMICs do not form links to the core 7TMIC cluster(s) (although all join a 7TMIC community structure by iteration 3 [**Figure 3B**]).

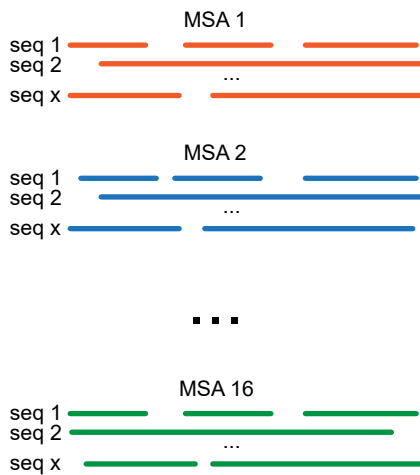
(B) Embedding-based conservation scores weakly but significantly correlate with column sequence identity from the A0A812K102-centered sequence alignment.

(C) Average embedding-based conservation scores for different subsets of 7TMICs, demonstrating that, while family-specific patterns exist, the conservation of anchor domain and pore regions is consistent. The TM and domain labels are derived from A0A812K102, as visualized in **Figure 3**.

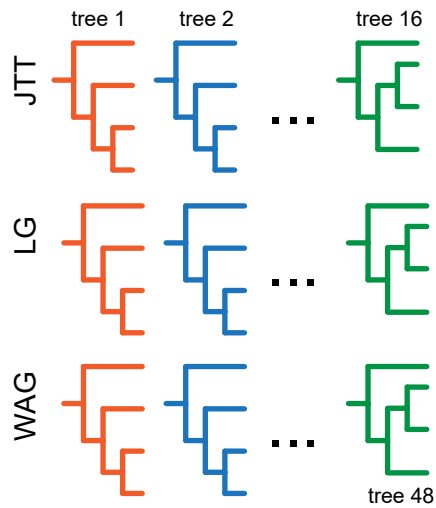
(D) Predicted tetramers for select 7TMICs. Top: top (presumed extracellular) and side views of the tetrameric arrangement of 7TMICs predicted by AlphaFold-Multimer, showing the formation of a hypothetical pore along TM7b, similar to *A. bakeri* Orco (far-left). Bottom: local Distance Difference Test (IDDT) scores (used to assess model confidence), plotted for each of the 5 replicate models generated. Each color represents a different replicate; vertical black lines separate each of the modelled subunits. Generally, the transmembrane-spanning alpha helices are the most confidently predicted, leading to the similar pattern of IDDT peaks and troughs across models.

**A**

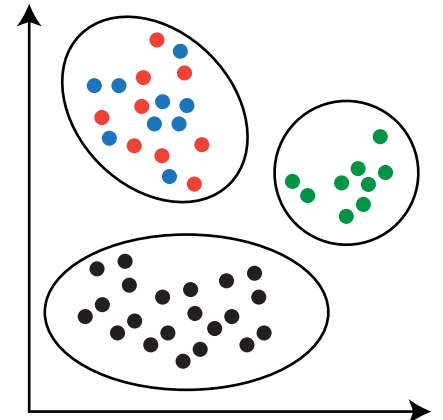
Generate ensemble of 16 MSAs by 4 guide tree and 4 HMM PRNG perturbations



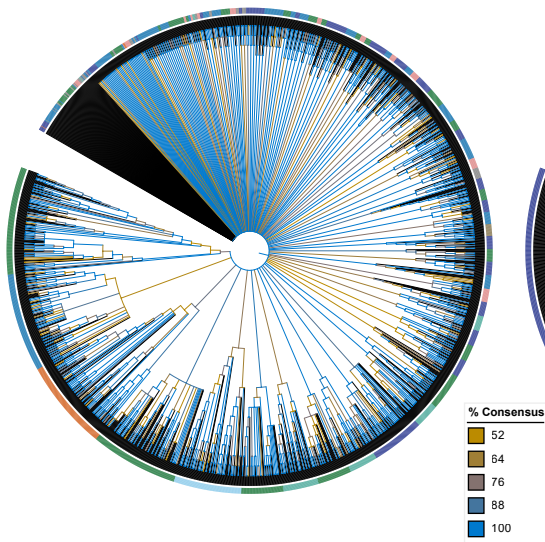
Generate forest of 48 phylogenetic trees using 3 substitution models



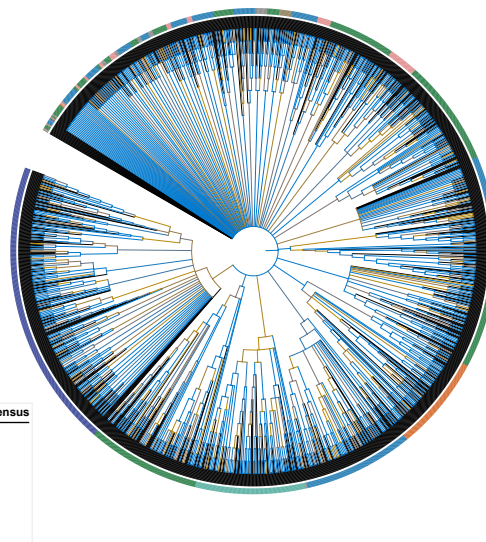
Calculate all-against-all Robinson-Foulds distances and visualize clusters of tree topologies by PCoA

**B**

all representative sequences

**C**

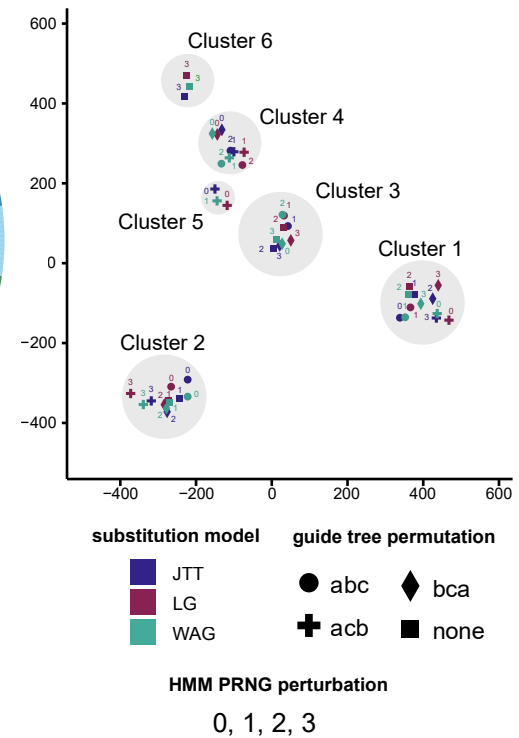
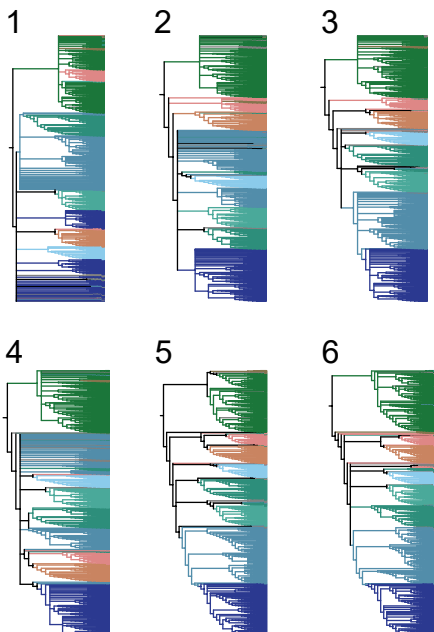
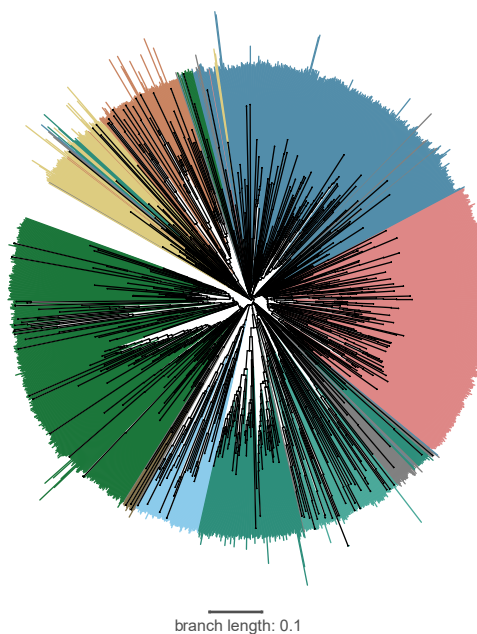
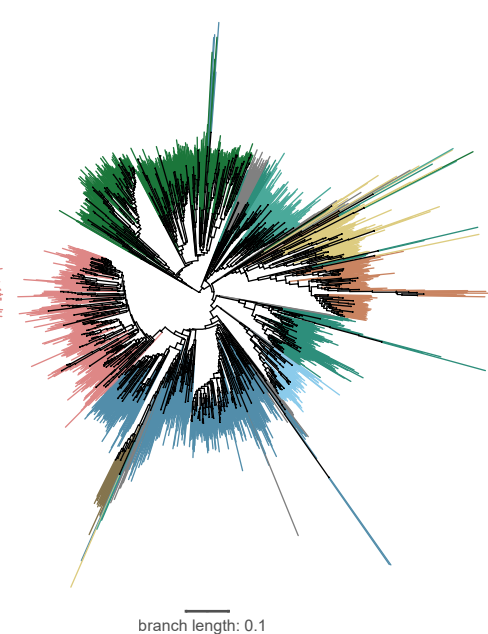
rogue taxa removed



Arch7TMIC Kineto7TMIC PHTF PHTF-like DUF3537 DUF3537-like  
 Bac7TMIC Alveolata GRL GrHz Gr/Grl/GRL Or incertae sedis

**D**

PCoA of phylogenetic forest

**E****F****G**

**Figure S3. Phylogenetic and tree space analysis, related to Figure 4.**

(A) Pipeline for sequence-based phylogenetic analysis. First, an ensemble of 16 multiple sequence alignments (MSAs) are made by perturbing the guide tree and the Hidden Markov model's pseudorandom number generator (HMM PRNG). Second, phylogenetic trees are generated for each of the MSAs, using 3 different amino acid substitution models, resulting in 48 trees. Finally, differences in the topology of the 48 trees are calculated by pairwise Robinson-Foulds distances; the resulting distance matrix is subsequently visualized in two dimensions by principal coordinate analysis (PCoA).

(B) Majority consensus tree for the 48 phylogenetic trees based on alignments of the representative 7TMIC sequences. 7TMIC clades/colors were assigned manually based on visual inspection of a CLANS-based clustering analysis; the color key for 7TMIC subfamilies is shown below panels (B-C). Branch colors indicate the percent consensus. There is essentially no clear consensus among these 48 initial trees; obviously monophyletic clades—such as insect Ors—are not reliably predicted, suggesting substantial alignment/phylogenetic errors (as expected for this highly divergent superfamily).

(C) Majority consensus tree for the 48 phylogenetic trees based on alignments of a 7TMIC dataset where rogue taxa (i.e. the most phylogenetically unstable leaves) have been removed. While there is still no greatly informative majority consensus topology, this analysis better recapitulates more obvious monophyletic clades, with higher branch consensus, indicating that errors have been minimized (but not eliminated, which we did not expect to occur at these levels of sequence dissimilarity).

(D) PCoA of Robinson-Foulds tree space for trees from (C). Trees form 6 topology clusters.

(E) Majority consensus trees for each of the 6 clusters, with colors matching (B) and (C). Five of these clusters agree that Kineto7TMICs branch proximally to prokaryotic 7TMICs, consistent with the hypothesis that kinetoplastids (and allies: Discoba) split early in eukaryotic evolution. Clusters 1 and 4 do not have majority consensus on deep 7TMIC branching. The remaining clusters suggest there are at least two Euk7TMIC families, termed Class-A and Class-B Euk7TMICs, but do not agree on the monophyly of Class-B Euk7TMICs. Clusters 4-6 suggest Class-B monophyly, while clusters 1-3 suggest that many proteins are basally branching (and thus, paraphyly). Given that structure-based phylogenetics suggest a monophyletic Class-B, this discordance may be the result of lingering long branch attraction or other errors resulting from the inclusion of rapidly evolved, horizontally-transferred and/or structurally-convergent proteins.

(F) Structural phylogeny derived from pairwise distances used the Foldseek 3Di structural alphabet, with colors matching the panels above. This tree is presented as rooted, but as in Figure 4, the true root is likely within the prokaryotic 7TMICs, at the location of the Last Universal Common Ancestor.

(G) Structural phylogeny derived from pairwise IDDT scores, with colors matching the panels above. As in (F), the true root is likely at location of the Last Universal Common Ancestor.



MSA Replicate	MSA Perturbations		All Sequences		No Rogue Taxa	
	Guide Tree	PRNG	Columns	Column Confidence	Columns	Column Confidence
1	abc	0	38154	0.315	36865	0.311
2	abc	1	38357	0.297	34833	0.295
3	abc	2	36970	0.314	35376	0.325
4	abc	3	37064	0.314	31247	0.304
5	acb	0	38375	0.31	35627	0.304
6	acb	1	40342	0.298	34754	0.302
7	acb	2	37080	0.317	34994	0.321
8	acb	3	35735	0.313	31345	0.305
9	bca	0	38391	0.315	35914	0.314
10	bca	1	39831	0.3	35813	0.294
11	bca	2	37647	0.316	34847	0.322
12	bca	3	36515	0.305	31406	0.303
13	none	0	38700	0.3	35842	0.309
14	none	1	40114	0.306	34808	0.293
15	none	2	37443	0.308	34932	0.309
16	none	3	37831	0.309	32689	0.298

**Table S1. Muscle5 multiple sequence alignment analysis, related to Figure 4.**

Column confidence is a measure of the reproducibility of each column, where 0 indicates the column is never found and 1 indicates it is found across all alignments. Dispersion is measured as the median dispersion of aligned letter pairs over the ensemble (D\_LP), and the median dispersion of columns over the ensemble (D\_Cols) (Robert Edgar, personal communication, 10 May 2023), where 0 is all the same and 1 is all different. Dispersion was extremely high. For the initial set of alignments: D\_LP=0.5836 D\_Cols=1.0000. After removal of rogue taxa: D\_LP=0.5855 D\_Cols=1.0000.