

# Methods for improving and characterizing gene families

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Gene families . . . . .	2
1.2	Improving and characterizing gene families . . . . .	2
1.3	Gene families in legumes . . . . .	3
1.4	Specific objectives . . . . .	3
<b>2</b>	<b>Objective 1: Develop pair classification based methods for detecting under-clustered gene families</b>	<b>4</b>
2.1	Methods . . . . .	4
2.1.1	Family searching and collecting closest non-family sequences . . . . .	4
2.1.2	Forming family and non-family pairs . . . . .	4
2.1.3	Model building and assessing classification performance . . . . .	4
2.1.4	Predicting missing sequences . . . . .	6
2.2	Preliminary results . . . . .	6
2.3	Significance . . . . .	8
2.4	Future work . . . . .	12
2.4.1	Merging small legume families . . . . .	12
2.4.2	Solving scalability issues . . . . .	12
<b>3</b>	<b>Objective 2: Develop tree based and domain composition based methods for detecting over-clustered families</b>	<b>13</b>
3.1	Tree based family scoring . . . . .	13
3.1.1	Methods . . . . .	13
3.1.2	Preliminary results . . . . .	13
3.1.3	Significance . . . . .	14
3.1.4	Future work . . . . .	14
3.2	Domain composition based family scoring . . . . .	15
3.2.1	Methods . . . . .	16
3.2.1.1	Domain composition based jaccard score . . . . .	16
3.2.1.2	Domain feature vector based cosine score . . . . .	17
3.2.2	Preliminary results . . . . .	17
3.2.3	Significance . . . . .	17
3.2.4	Future work . . . . .	23
3.3	Species composition based family scoring . . . . .	23
3.3.1	Methods . . . . .	23
3.3.2	Preliminary results . . . . .	23
3.3.3	Significance . . . . .	24
3.3.4	Future work . . . . .	24

<b>4</b>	<b>Objective 3: Develop a workflow for building consensus families from multiple family sets</b>	<b>26</b>
4.1	Methods . . . . .	26
4.1.1	Building HMM sets for families . . . . .	26
4.1.2	Classifying common set of proteomes into families using HMM sets . . . . .	26
4.1.3	Calculating family correspondence and family overlaps between the sets . . . . .	26
4.2	Preliminary results . . . . .	26
4.3	Future work . . . . .	27
<b>5</b>	<b>Objective 4: Biological characterization of the legume families</b>	<b>29</b>
5.1	Methods . . . . .	29
5.2	Preliminary results . . . . .	30
5.3	Future work . . . . .	30

## Abstract

Gene families are groups of genes that have evolved from single common ancestral gene present in the common ancestor of a given set of species. Studying evolution of genes in the families helps in inferring evolution of species and in detecting evolutionarily/functionally equivalent regions between multiple different genomes. Current gene family clustering methods can produce non-optimal families that could be under-clustered (missing sequences) or over-clustered (containing sequences from more than one family). This project presents methods for analyzing gene families - for over- and under-clustering and by other characteristics, after they have been built using existing family building methods. The proposed methods include (1) a pair classification based approach for detecting and correcting under-clustered families; (2) tree-based methods that assess families based on expected phylogenetic topology; (3) domain composition methods that assess families based on domains detected in the sequences in a family; (4) species composition methods that assess families based on expected species in a family, relative to a model of gene duplications and speciations and (5) family sets comparison workflow for obtaining consensus families from multiple family sets built using different methods. Finally, the high-quality gene families from legume species, obtained using the proposed methods, are first divided into different groups based on their sizes and each group is characterized using Gene Ontology enrichment analysis.

# 1 Introduction

## 1.1 Gene families

Gene families or orthologous groups - with respect to a given set of species - are clusters of genes that have originated from a single ancestral gene present in the common ancestor of the given set of species. Therefore, gene families contain genes that have diverged due to speciation or duplication after the divergence node of earliest diverging species. Accordingly, genes within the families are either orthologs or in-paralogs (duplication after the common ancestral node) (Fitch, 1970, 2000; Sonnhammer and Koonin, 2002). Gene families are mainly used to infer orthology and paralogy among genes from the given set of species.

Clustering genes into families helps in capturing evolutionary information from multiple species for studying the divergence or conservation of particular biological functions, with respect to given set of species. Gene families can also be used for identification of candidates drug/vaccine development in case of parasites like *Plasmodium falciparum* which causes malaria in humans (Gardner et al., 2002; Kissinger et al., 2002; Whetzel et al., 2005). Gene families have also been used for annotation of newly sequenced genomes and for anti-bacterial drug development by cross-referencing functional information from multiple species (Tatusov et al., 1997; Galperin and Koonin, 1999; Natale et al., 2000b,a; Forterre, 2002)

## 1.2 Improving and characterizing gene families

The clusters of Orthologous Groups (COG) database is one of the earliest widely-used collections of gene families (Tatusov et al., 2000, 2001). The orthologous groups in the COG database were identified based on reciprocal best hits from all-against-all BLAST searches using complete proteomes from five major lineages - Gram-negative bacteria, Gram-positive bacteria, Cyanobacteria, Archea, and Eukarya. Existing family

building algorithms can be divided into 2 main categories, viz: alignment-based clustering approaches and tree based approaches. The alignment-based clustering approaches either operate on individual pairs of species and then expand to multiple species - for example MultiParanoid (Inparanoid), OMA (Alexeyenko et al., 2006; Remm et al., 2001; Roth et al., 2008; Altenhoff et al., 2010) - or attempt to infer families directly from multiple species with the common ancestor of given set of species as a reference; OrthoMCL, OrthoFinder (Li et al., 2003; Emms and Kelly, 2015). Tree-based methods compare the gene trees with the species trees using tree reconciliation methods for building gene families (Li et al., 2006; Van der Heijden et al., 2007; Wapinski et al., 2007; Hubbard et al., 2006; Dehal and Boore, 2006).

While many of the families produced by existing methods may be "correct", many non-optimal families could be produced due to under-clustering or over-clustering. Under-clustered families are those families that may have been fragmented by the clustering method (missing sequences). Over-clustered families are those families that may have been produced due to merging of two or more closely related families. Also, a significant number of sequences can remain un-clustered and could be wrongly considered as singletons or orphan genes.

Therefore, it is important to analyze and assess the accuracy of families constructed using any method. Here, I propose a suite of methods and tools for improving under-clustered/over-clustered gene families, and for characterizing true and consistent families.

### 1.3 Gene families in legumes

The legume family (Fabaceae) consists of 750 genera and 19,500 species of flowering plants, making the third largest family of flowering plants (Lewis, 2005). The family is broadly divided into four major lineages viz. Papilionoideae, Mimosoideae-Cassiinae-Caesalpinieae (MCC), Detarieae and Cercideae, out of which the Papilionoideae is the largest. The Papilionoid sub-family includes species like *Arachis hypogaea* (cultivated peanut), *Cajanus cajan* (pigeonpea), *Cicer arietinum* (chickpea), *Glycine max* (soybean), *Lotus japonicus* (Lotus), *Lupinus angustifolius* (narrow-leaved lupin), *Medicago truncatula* (barrel medic), *Phaseolus vulgaris* (common bean), *Trifolium pratense* (red clover), *Vigna angularis* (adzuki bean), *Vigna radiata* (mungbean) and *Vigna unguiculata* (Cowpea).

An ancient Whole Genome Duplication (WGD) has occurred in common ancestor of the Papilionoid sub-family around 55 Ma. (Blanc and Wolfe, 2004; Schlueter et al., 2004; Pfeil et al., 2005; Cannon et al., 2006; Bertoli et al., 2009; Cannon et al., 2014). In addition, the lineages such as *Glycine* and *Lupin* have also undergone independent lineage specific WGDs (Cannon et al., 2014; Kroc et al., 2014).

Approximately 17k-18k gene families were obtained using a Ks-based (synonymous substitutions per site) family building procedure from the Papilionoid species proteomes. Additionally, for comparison, a different set of 20k families was also obtained from Papilionoid species using the OrthoFinder method (Emms and Kelly, 2015).

### 1.4 Specific objectives

1. Develop pair classification based methods for detecting and improving under-clustered gene families.
2. Develop tree based, domain composition and species composition based methods for detecting and improving over-clustered gene families.
3. Develop a workflow for building consensus families from multiple family sets.
4. Biological characterization of gene families for a set of species in the legume plant family.

## 2 Objective 1: Develop pair classification based methods for detecting under-clustered gene families

This scoring method leverages the basic evolutionary property of gene families - all the sequences in the family must have diverged at/after the separation of earliest diverging species and none of the sequences outside the family have diverged in the family. Accordingly, for a given family, this method attempts to separate pairs from the family, from pairs where one sequence is in the family and other is outside the family. The separation statistics tell how well a given family of sequences can be circumscribed from closest non-family sequences. This can provide an indication of family completeness.

### 2.1 Methods

#### 2.1.1 Family searching and collecting closest non-family sequences

All the sequences from any given family are searched against the combined database of all proteomes from all the species under study (the sequence space). With each family sequence as query, all the hits that match the query better than the worst matching family hit (other than the query) are collected into a list. A combined list of hits, using each family sequence as a query, is obtained which is expected to contain both family and non-family sequences (Figure 1).

The phmmer program from HMMER package (Eddy, 1992) is used for searching families against the sequence space database. The ranking of hits for finding the worst matching family hit is the same ranking given by the phmmer search results.

#### 2.1.2 Forming family and non-family pairs

The combined list of family and non-family sequences obtained in the previous step is used to form family and non-family pairs (Figure 2). The family pairs are those that form exclusively between the original family members and the non-family pairs are those that form between the family and non-family members. The family pairs are labeled as positive pairs and the non-family pairs were labeled as the negative pairs.

#### 2.1.3 Model building and assessing classification performance

HMM based pair classification models are built to classify the family and non-family pairs, and 10 iterations of repeated test-train split strategy are used to assess the classification performance (Figure 3). For each iteration, the set of family+non-family pairs are randomly split into training and test split. The pair classification model is trained on the training split and tested on the unseen test split. To train the HMM model on positive pairs in the training split, consensus sequences of the positive pairs are used. To test the model on pairs in the test split, individual sequences of the pairs are first aligned to the trained HMM model and alignment scores for both the sequences are obtained. The test pair is predicted as a positive pair if both the alignment scores are greater than or equal to a specified score cutoff.

The Precision-Recall (PR) curve for the positive class is obtained using the combined predictions from all the 10 test-train split iterations. The test pair alignment scores from all the iterations are consolidated (Figure 4) and used for calculating precision ( $TP/(TP+FP)$ ), recall ( $TP/(TP+FN)$ ) and F1-score (harmonic mean:  $2 \cdot (\text{precision} \cdot \text{recall}) / (\text{precision} + \text{recall})$ ) values for specified score cutoffs, where TP, FP and FN are the number True Positive, False Positive and False Negative pairs, respectively. Different precision and recall values are obtained for a range of score cutoffs starting from most stringent to least stringent. The precision values are plotted against the recall values to obtain the PR-curve (Davis and Goadrich, 2006). The area under the PR-curve (PR-AUC) is calculated using the trapezoidal rule. Other classification metrics such as the precision and recall values observed at the best F1-score and the score cutoff which gives the best F1-score are also reported. The lowest score observed for the positive class is also reported as the lowest score cutoff with respect to the given family.

An example of PR-curve plot for a hypothetical gene family is shown in Figure 5. Each point on the curve corresponds to an alignment score cutoff with score cutoffs decreasing from left to right. The highest score cutoffs on the left produce classifications with highest precision (low FP) but lowest recall (high FN). Conversely, lowest score cutoffs on the right produce classifications with lowest precision but highest recall.

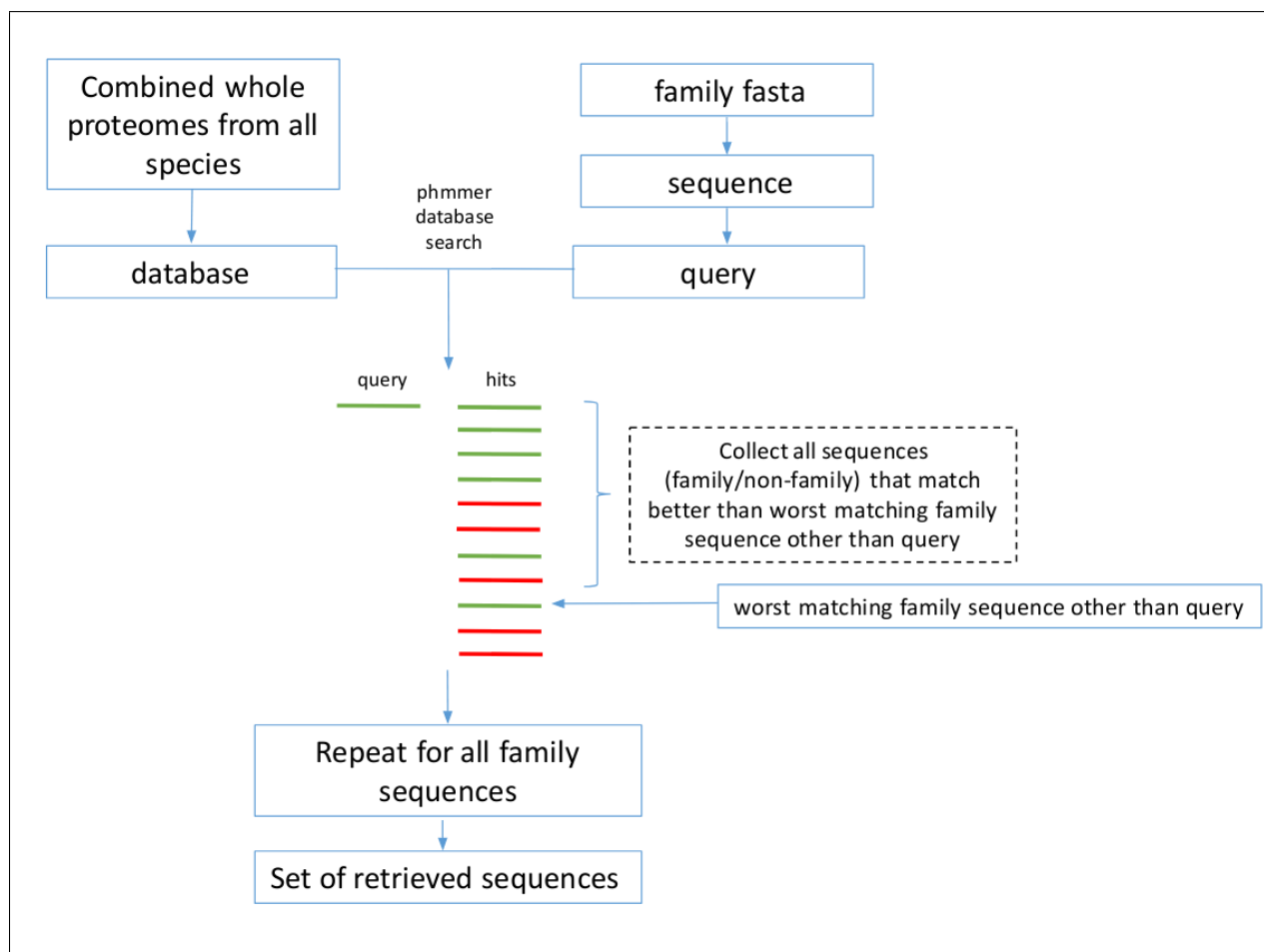


Figure 1: Collecting closest non-family sequences

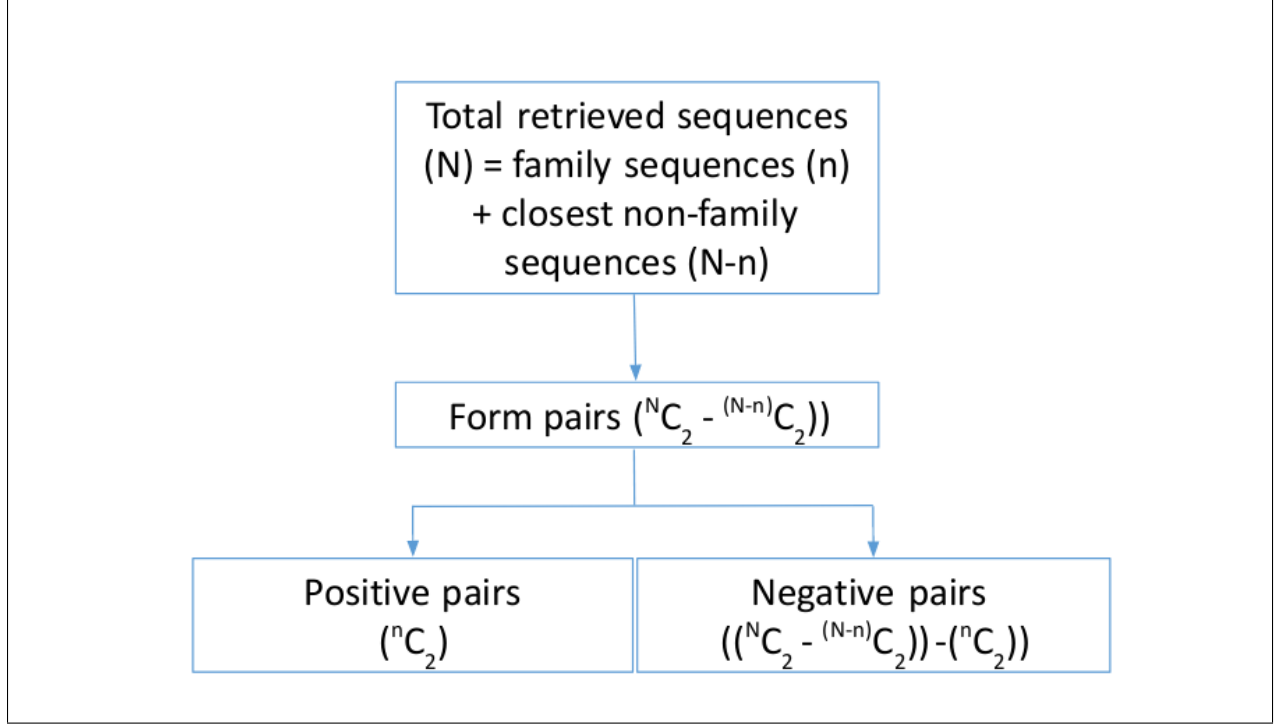


Figure 2: Forming positive and negative pairs from family and non-family sequences

An F1-score can be calculated for each point (precision, recall) on the curve as the harmonic mean of precision and the corresponding recall value. The point on the curve with the highest F1-score is the point where optimal values of precision and recall exist. This point represents the best classification performance for the family and the corresponding alignment score is the score cutoff that gives the best classification performance for the family.

#### 2.1.4 Predicting missing sequences

Candidate missing sequences for every family are predicted using the non-family pairs. All the non-family pairs where the scores for both the sequences of a pair are greater than the best/lowest score cutoff are re-classified as positive pairs. Unique sequences within these reclassified positive pairs are predicted and reported as candidate missing sequences for the family.

## 2.2 Preliminary results

The pair classification based scoring was tested on 4796 yeast families from the Yeast Gene Order Browser (YGOB) database (Byrne and Wolfe, 2005). Since the YGOB families are built through manual curation using synteny based evidence, they were assumed to be "true". To check if the quality assessment workflow is assigning high quality scores to all the true families, the distribution of the PR-AUCs for all the 4796 yeast families was obtained (Figure 6). As expected the distribution of PR-AUCs is highly skewed towards PR-AUC value = 1.0 with 92% of family classifiers having values  $\geq 0.75$ . This shows that the method correctly recognizes good quality families and assigns high quality scores to them.

Even though the quality assessment method is proven to work on good/true families it is important to study the behaviour of the method on non-optimal families. To check the behaviour of the method on incomplete families, the true yeast families were artificially manipulated so that all the families are randomly missing 20% of family sequences.

The histogram in Figure 7 shows the distribution of PR-AUC values for 4796 artificially manipulated families where every family is missing 20% of their sequences. A plot (Figure 8) of true families PR-

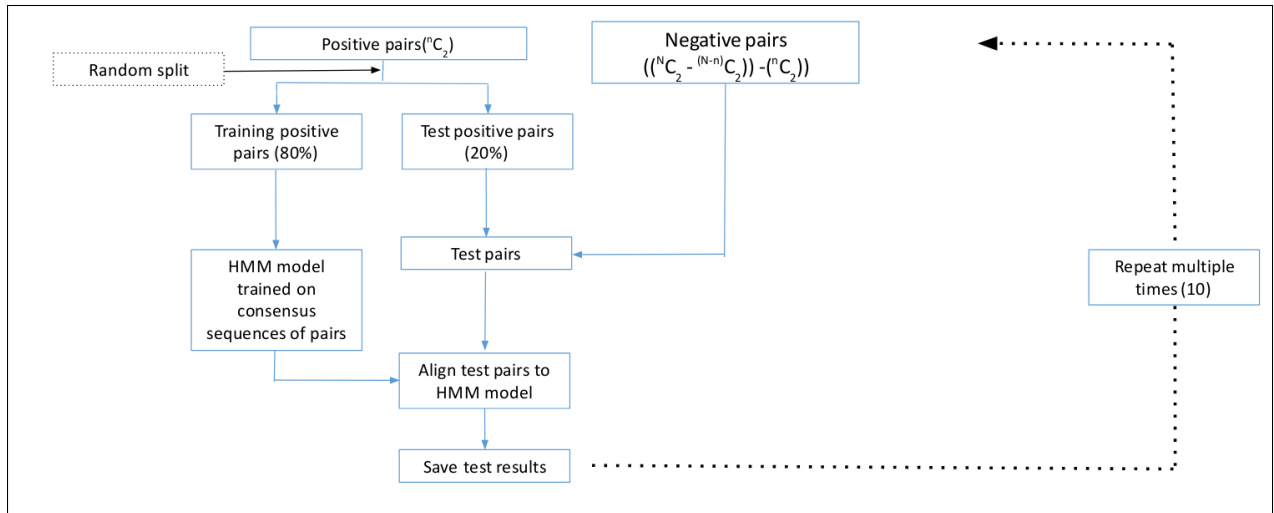


Figure 3: Training pair classification models for family

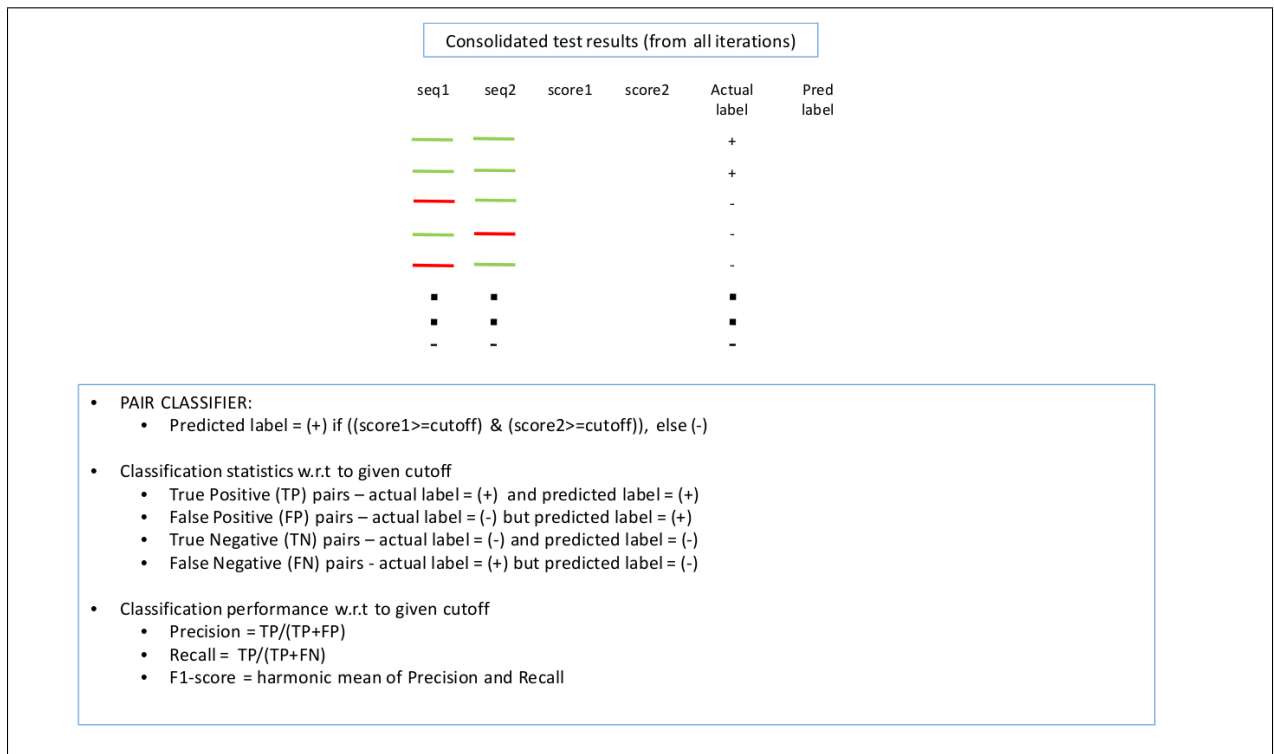


Figure 4: Evaluating pair classification model

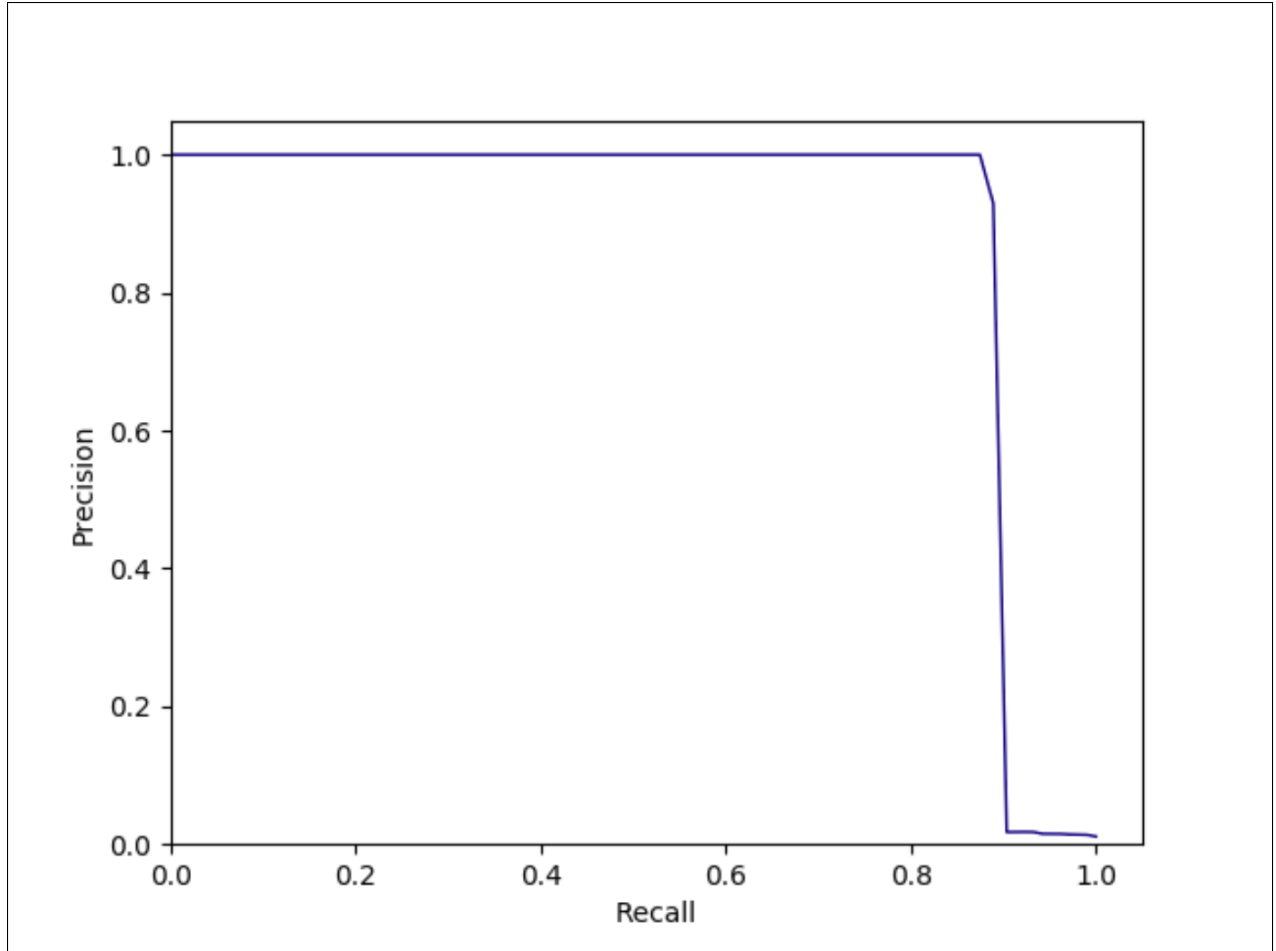


Figure 5: Example of PR-curve obtained during training

AUCs versus PR-AUCs after removing family sequences shows the drop in the PR-AUCs after removing family sequences from true families. The distribution indicates that the overall quality of the families drop significantly, which shows ability of method to detect incomplete families. Specifically, for 3971 out of 4796 families (83%), the PR-AUC dropped significantly ( $\text{PR-AUC} \geq 0.9$  for true families and  $\text{PR-AUC} < 0.75$  after removing family sequences)

For each manipulated family, the missing sequences were predicted back using lowest score cutoff obtained during training. The predicted missing sequences were compared to the true missing sequences and, precision and recall values were calculated to study the accuracy of predicted missing sequences. Precision was calculated as  $(\text{TP}/(\text{TP}+\text{FP}))$  and recall was calculated as  $(\text{TP}/(\text{TP}+\text{FN}))$  where True Positive (TP) are those predicted missing sequences that were truly missing from the family, False Positive (FP) are those that were predicted as missing but do not actually belong the family and False Negative (FN) are those that are truly missing but were not predicted as missing by the method. Out of 4796, the prediction performance for missing sequences was good (precision  $\geq 0.75$  and recall  $\geq 0.75$ ) for 3760 (78.4%) families.

### 2.3 Significance

This family scoring method assumes the given family to be under-clustered - the majority of sequences have diverged into the family - and attempts to build an HMM-based classifier that classifies pairs of sequences formed within the family from the pairs formed between family and the closest non-family sequences. The classification performance is mainly given by the PR-AUC value which shows how well the given family is separable from the closest families. This scoring method can be used in two ways; to characterize true



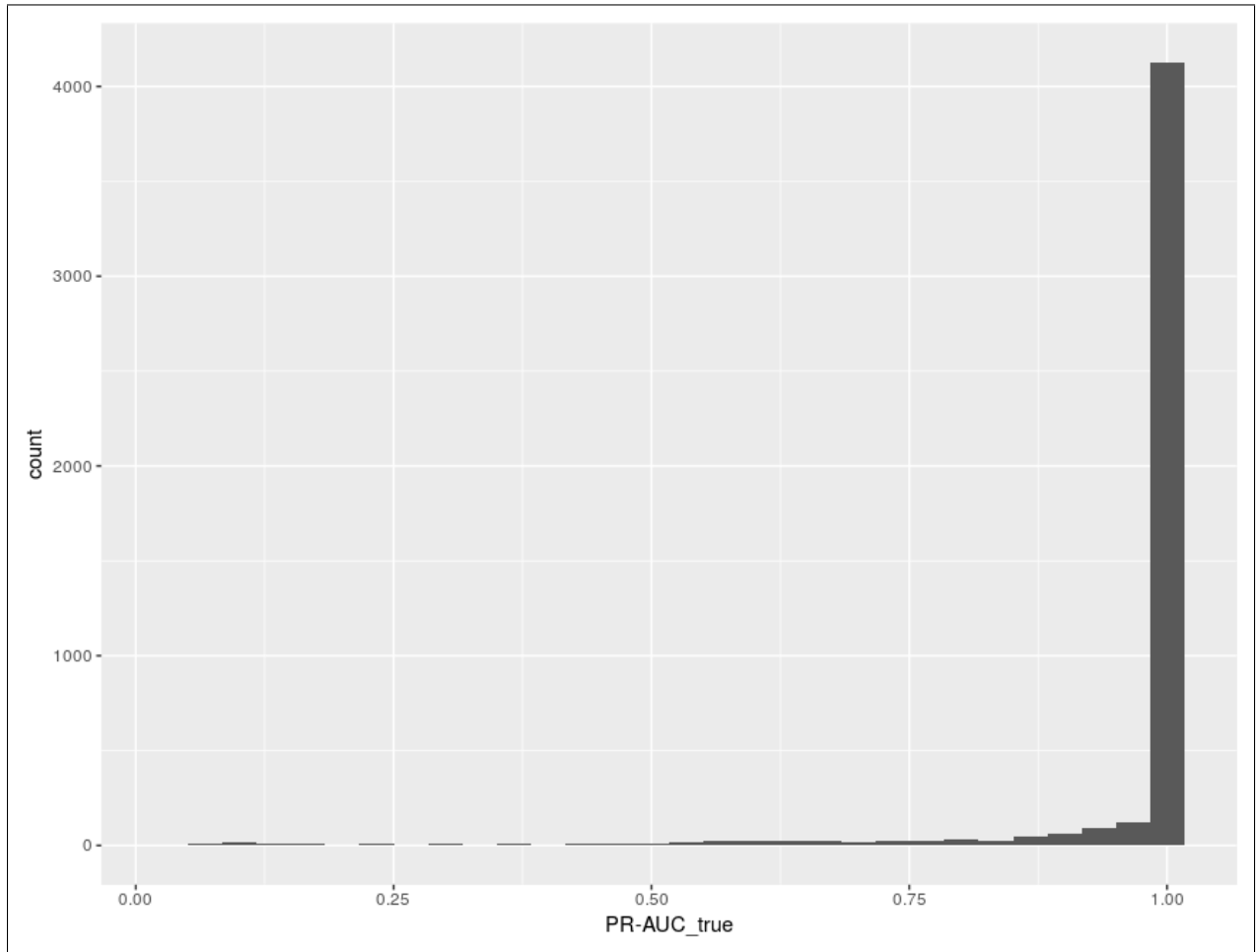


Figure 6: Distribution of PR-AUC values for true YGOB families

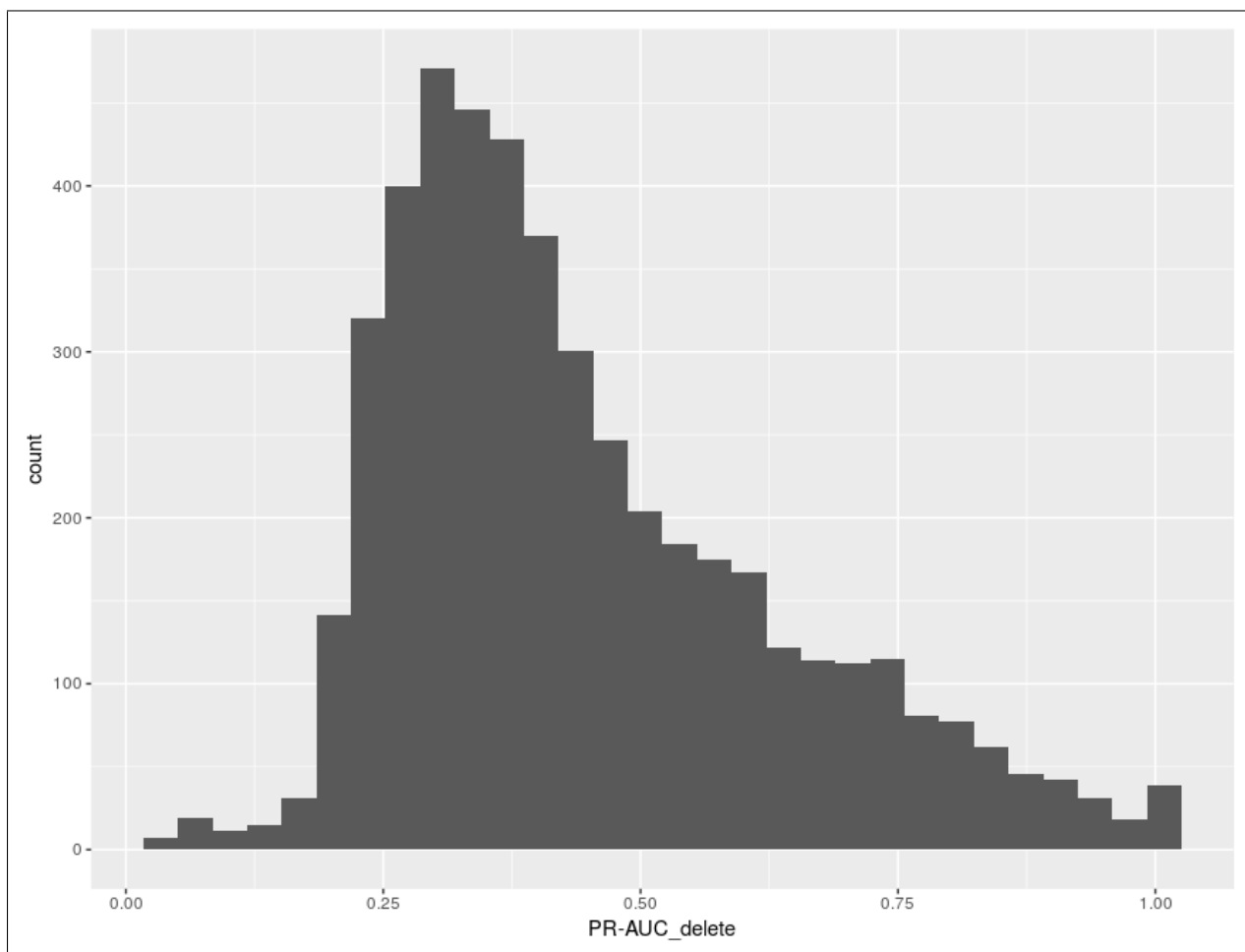


Figure 7: Distribution of PR-AUC values for families with 20% of sequences deleted

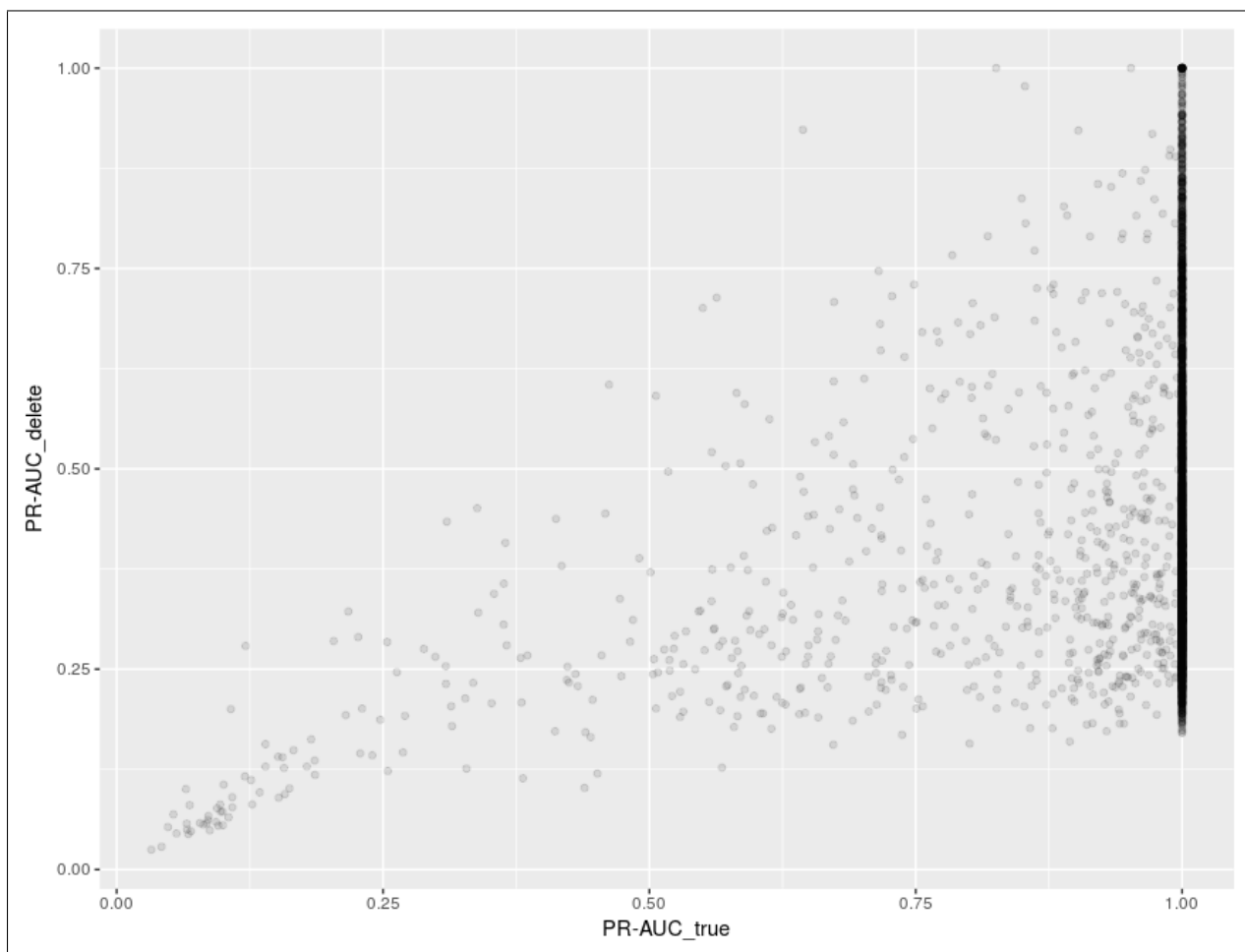


Figure 8: Plot of PR-AUCs obtained from true families versus PR-AUCs after deleting 20% of sequences

families and to detect missing sequences for under-clustered families. For characterizing true families, the PR-curve and the PR-AUC can be used to see how well-behaved the family is in terms of its evolution. The best score cutoff obtained during training can be used to add sequences to the families from newly sequenced species. For under-clustered families, the PR-AUC is expected to be low as the precision value drops due to greater number of FP pairs (sequences trying to get in the family). The missing family sequences can be predicted using the best or the lowest score cutoffs depending upon preference of high precision or high recall for prediction. If precision is preferred over recall then the best score cutoff can be used and if recall is preferred over precision then the lowest score cutoff can be used for predicting missing sequences.

## **2.4 Future work**

### **2.4.1 Merging small legume families**

This scoring method will be used to inspect and merge small legume families into larger families. In the Ks-based legume families, there are approximately 2000 families that may be under-clustered ( $\leq 10$  family size). This will help in reducing the number of under-clustered families and generating a set of high-quality legume families.

### **2.4.2 Solving scalability issues**

The pair classification based method is not fast enough to be efficiently applied on all legume the families. There are mainly two steps in the workflow that can be optimized to speed-up the method viz. the initial phmmer based database search and selecting the closest non-family sequences for forming negative pairs. Currently, every given family is searched against the entire set of all proteomes under study. The search space can be reduced by only searching within the neighborhood sequences/families. Tools like HHsearch (Söding, 2004) can be used to detect neighborhood families and limit the search space for every family.

Similarly, when selecting the closest non-family sequences, only those non-family sequences that are detected by more than one family sequences (or some proportion of family sequences) can be selected. Currently, a union of non-family sequences detected by every family sequence are selected for forming negative pairs.

### 3 Objective 2: Develop tree based and domain composition based methods for detecting over-clustered families

#### 3.1 Tree based family scoring

Family trees are largely expected to follow the pattern of evolution from the species tree. This correspondence between gene trees and species trees has been used previously in multiple methods like TreeFam, LOFT, SYNERGY, ENSEMBLE 2007, PhIGs (Li et al., 2006; Van der Heijden et al., 2007; Wapinski et al., 2007; Hubbard et al., 2006; Dehal and Boore, 2006) for building gene families and detection of orthologs and paralogs from the family trees. Animal gene families in TreeFam are the group of genes evolved from a single gene in last common of all animals. Genes in one family are identified on the basis that they are phylogenetically separated by one or more genes from the outgroup species (*S. cerevisiae*, *S. pombe*. or *A. thaliana*)

Since gene families contain sequences that have diverged at or after the divergence of the earliest diverging species under study, genes within families have also diverged after the divergence of outgroup species. This evolutionary property of gene families can be used to score family trees by calculating the proportion of ingroup sequence pairs that appear to have diverged after the divergence of outgroup sequences in the corresponding phylogenetic trees of the families.

##### 3.1.1 Methods

For a given rooted family tree, each pair of ingroup sequences, found within the tree, is labeled as True Positive (TP) or False Positive (FP) depending upon whether the pair appears to have diverged after or before the divergence of one or more outgroup sequences. A precision score for the family is calculated as  $TP/(TP+FP)$  which gives the proportion of ingroup pairs diverging after the outgroup separation in the family tree.

To check the divergence of any ingroup sequence pair in the tree, the Most Recent Common Ancestor (MRCA) of the pair is obtained. Then, all sequences corresponding to leaf nodes under this MRCA are collected. If this set of sequences contains one or more outgroup sequences, the corresponding ingroup pair is labelled as FP, else it is labeled as TP. The FP label for any ingroup pair indicates that there is atleast one outgroup sequence that has diverged after the divergence of the ingroup pair and the corresponding sequences of the pair have been wrongly clustered into one family.

##### 3.1.2 Preliminary results

This scoring method was applied on 18543 legume families. The Figure 9 shows the distribution of scores across the legume families. Most of the families have higher precision scores ( $\geq 0.75$ ).

To study the relationship between the tree based pair precision scores and family size, plot (Figure 10) of precision scores versus number of ingroup sequences was obtained. As expected, most of families containing 15-45 ingroup sequences have high precision scores. A considerable number of families, containing approximately 40 ingroup sequences, have tree precision scores around 0.5. These are the families where approximately half of ingroup pairs diverge after the divergence of one or more outgroup sequences. These families are prime candidates of over-clustered families. Also, large families with very low precision scores (majority of ingroup pairs have diverged before the outgroup divergence) may indicate problems with tree rooting.

Examples of family trees with typical values of tree based pair precision scores are shown below. The outgroup sequences in the trees are shown in grey. Figure 11 shows family tree with low pair precision score (0.095). For this family, the majority of ingroup sequences do not form a monophyletic clade. Figure 12 shows family tree with pair precision score = 0.507. This family has two monophyletic ingroup clades that could potentially be two ingroup families clustered together by the family building algorithm. This shows that the pair precision score can detect such over-clustered families. Finally, Figure 13 shows a family with high pair precision score (0.876). Majority of ingroup sequences form a monophyletic clade with two sequences diverging outside the family.

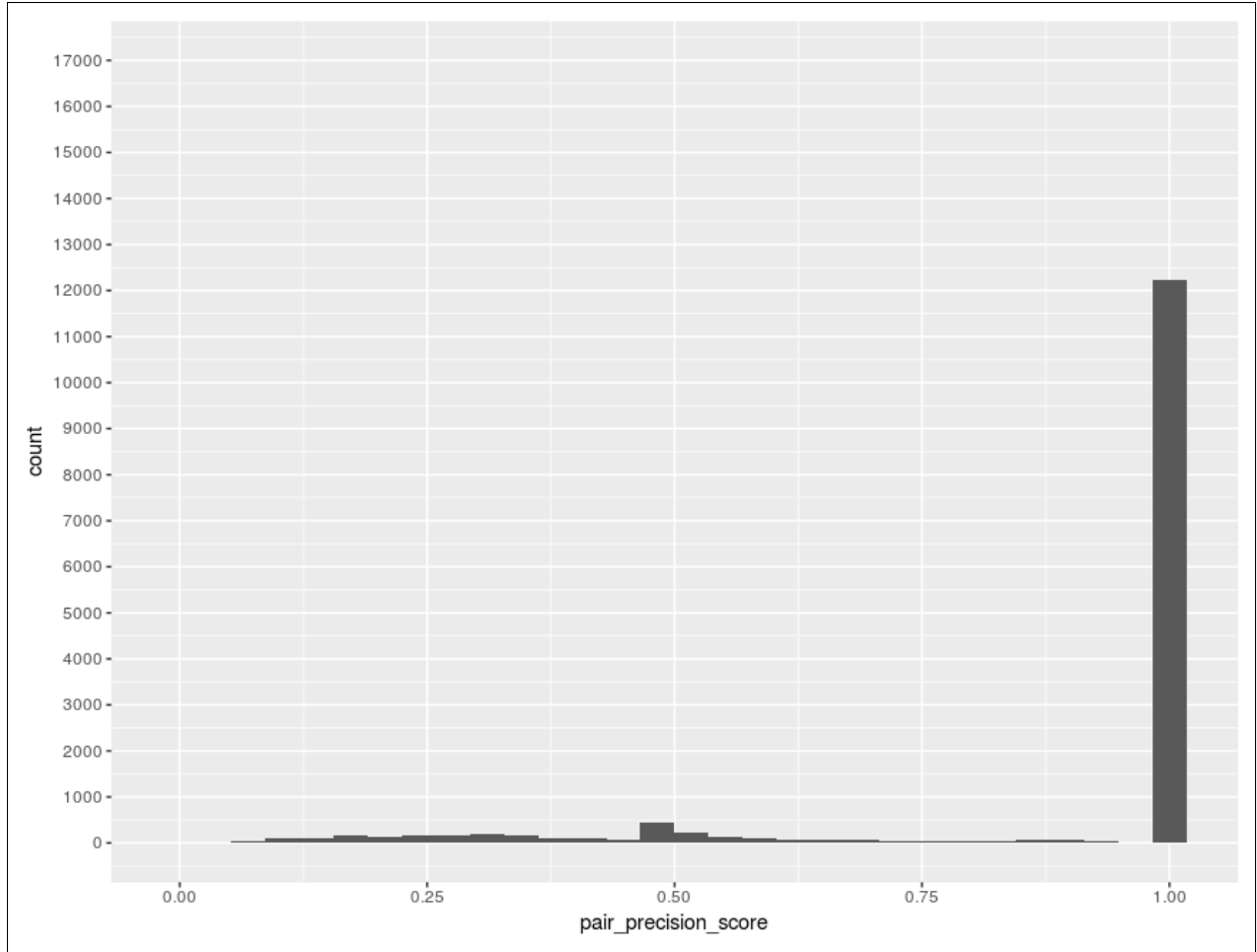


Figure 9: Distribution of tree based precision scores for Ks-based families

### 3.1.3 Significance

This simple tree-based scoring method uses basic evolutionary properties of gene families to detect if all the sequences clustered in the family truly belong in the family, using the rooted evolutionary family tree. As seen from the preliminary results, this scoring method can be used to detect over-clustered families and family trees with rooting problems.

### 3.1.4 Future work

Over-clustered legume families detected using this method will be corrected to obtain a set of high-quality legume families. Also, large families with low precision scores will be re-rooted to maximize the precision scores.

The merged families obtained using the pair classification based method will also be analyzed using the tree based scoring to resolve any over-clustered families produced due to merging.

An R/Python tool can be developed where users can analyze and score family trees with respect to a given set of outgroups. The tool will also be able to directly extract ingroup clades from over-clustered families.

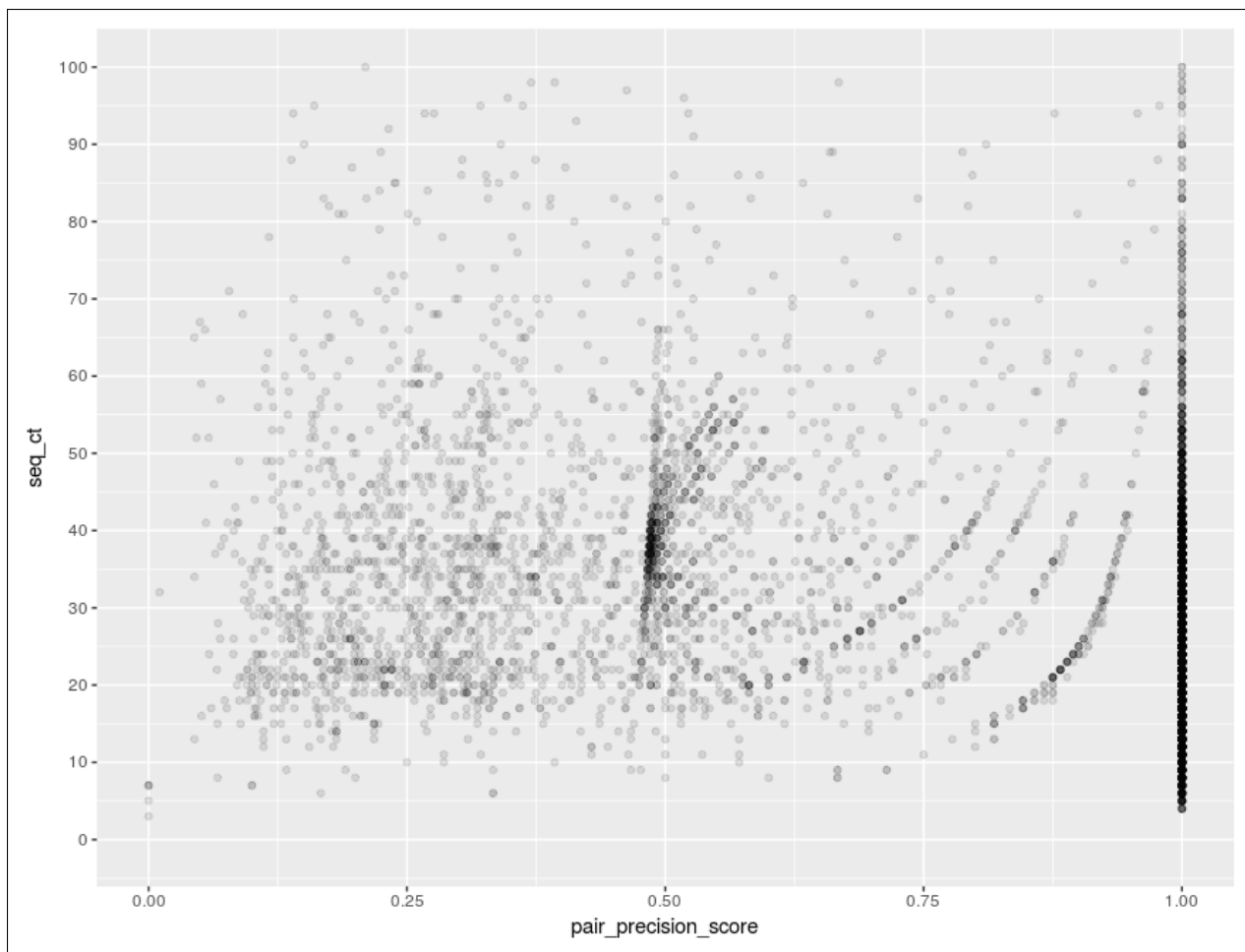


Figure 10: Plot of tree based pair precision scores versus ingroup sequence count for Ks-based families

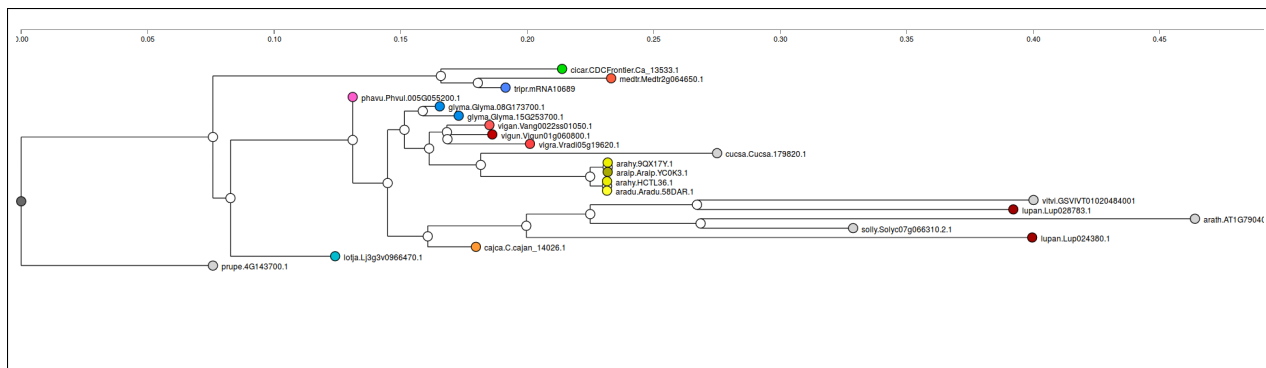


Figure 11: Example family tree with tree based pair precision score = 0.095

### 3.2 Domain composition based family scoring

Protein domains are sections of protein sequences that can fold and function independently. Consequently, multidomain proteins can evolve in a modular fashion through domain deletions/insertions/duplication in addition to sequence-based evolution. Databases like Pfam, PROSITE, SMART, CDD (Finn et al., 2007; Falquet et al., 2002; Letunic et al., 2011; Marchler-Bauer et al., 2010) define and store domains detected in large number of protein sequences. Domain compositions/content have been previously used to detect

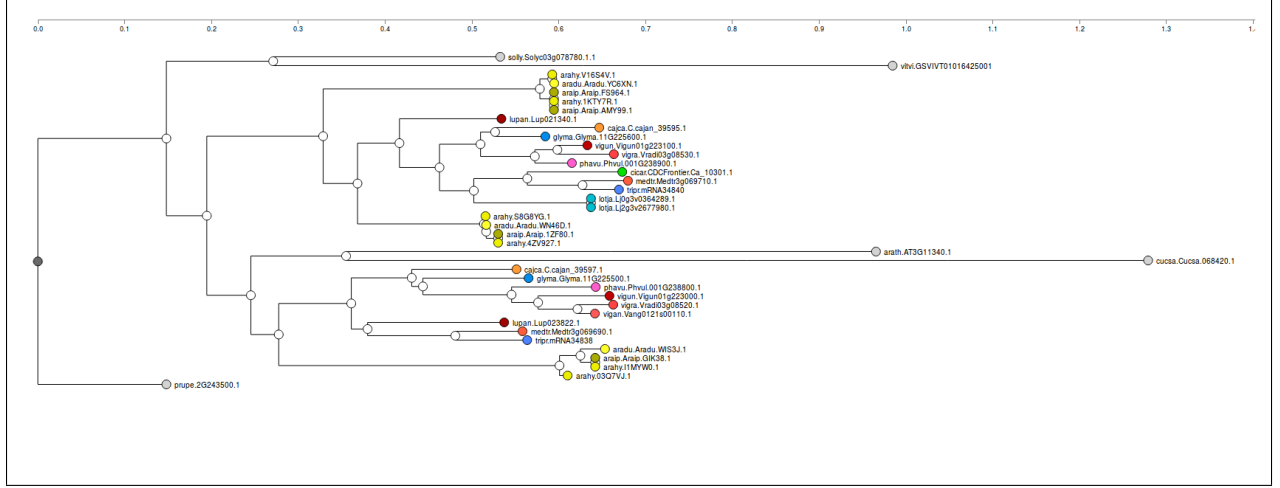


Figure 12: Example family tree with tree based pair precision score = 0.507

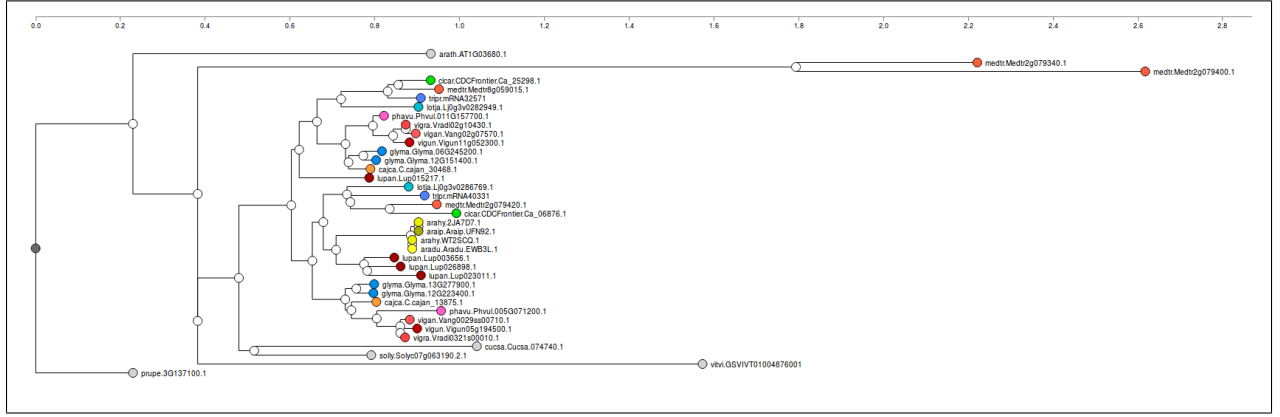


Figure 13: Example family tree with tree based pair precision score = 0.879

sequence homology with high accuracies. (Song et al., 2007; Bitard-Feildel et al., 2015). These were well separated gene families with distinctly different domain compositions (E.g. families separated at metazoan ancestor)

Closely related sequences from the same family can be expected to have fairly similar domain types and domain content. Accordingly, this method assigns scores to given families based on number of domains shared between the sequences of the family. This scoring scheme will assign higher scores to family where most of the sequences have similar domain compositions. Good and well-conserved families are expected to have higher domain composition scores. Conversely, over-clustered families containing diverse sequences are expected to have lower scores due to less conservation of domain compositions across all the sequences.

### 3.2.1 Methods

#### 3.2.1.1 Domain composition based jaccard score

For a given family, Pfam-A domains are assigned to all the sequences in the family using the pfamscan program (Mistry et al., 2007). The domains are assigned to sequences in order of their start coordinates in the sequences. Then, for each pair of sequences in the family, A domain composition based Jaccard score is calculated using the function:  $(n_1 \cup n_2) / (n_1 \cap n_2)$ , where  $n_1$  and  $n_2$  are the sets unique domains from sequence 1 and 2 of the pair, respectively. Jaccard score for the family is calculated as mean of all pairwise Jaccard scores from the family.



### 3.2.1.2 Domain feature vector based cosine score

Here, the domain compositions of sequences within a given family are compared using their domain feature vectors. This is a new scoring method for comparing domain compositions of two sequences using their domain feature vectors. For example, consider two sequences X and Y from the same family having following domain order along the sequence, X: {A,B,B,C} and Y: {A,A,B,D}. The duplicate domains in both the sequences are assigned unique ids to distinguish them from each other, X: {A, B, B-2, C} and Y: {A, A-2, B, D}. The domain content universe for both the sequences is (A, A-2, B, B-2, C, D). Accordingly, the domain feature vector for both sequences is X:  $(x_1, 0, x_2, x_3, x_4, 0)$  and Y:  $(y_1, y_2, y_3, 0, 0, y_4)$ , where  $x_i$  and  $y_i$  are the alignment scores for the corresponding domain HMMs aligning against the sequences X and Y, respectively. Cosine similarities are calculated between all pairs of sequences using their domain feature vectors and the domain composition score for the family is calculated as the mean of all pairwise cosine scores.

Pfam clusters closely related domains in Pfam-A into clans. For sequences within the same family, closely related domains could be detected interchangeably, for same positions, in closely related family sequences. To avoid this, clan-ids of the domains are used instead of domain-ids for calculating the domain feature vector based cosine scores.

### 3.2.2 Preliminary results

The domain composition based Jaccard and cosine scoring were applied on the existing legume families to study the relationships between different scores and the ingroup sequence count. The score distributions for domain-based Jaccard scores can be seen from Figure 14.

To study the relationship between the domain Jaccard score and family size, the plot of jaccard score versus number of ingroup sequences in the family was obtained (Figure 15). The Jaccard scores are high for the majority of families in which the number of ingroup sequences is between 15-45.

The distribution of domain feature based cosine scores was also obtained (Figure 16). Most of the families have high scores indicating that for a majority of the families, the domain compositions are fairly conserved.

The plot (Figure 17) showing the relationship between the cosine score and family size was also obtained to study the relationship between cosine score and number of ingroup sequences in the family. Here, the scores correlate well (better than Jaccard scores) with family sizes, with higher scores for families with number of ingroup sequences between 15-45.

A plot (Figure 18) showing the relationship between the tree based pair precision scores and domain cosine scores was also obtained. The majority of families with high tree based precision scores also have high cosine scores. Specifically, for 11899 families, out of 15428 families, (77.1%) both, the tree based pair precision scores and domain composition based cosine scores are  $\geq 0.75$ . In addition, two types of families can be seen from the plot; families with high tree based precision scores and lower domain scores and vice versa. The former are those that seem to evolve rapidly through domain insertion/deletion. The latter are those that seem to evolve slowly conserving domain compositions over longer evolutionary periods which could be the reason for the over-clustering.

Both the domain based family scoring methods have been developed into a Python tool that accepts the Pfam database and a directory containing family fasta files, and assigns domain composition scores to families.

### 3.2.3 Significance

Here we explore different domain based family scores mainly to detect over-clustering in families. Assuming that the species under study are closely related, the domain compositions for all the sequences in the family are expected to be largely conserved. The domain feature vector based cosine score seems to correlate better with family size and the tree based precision score as compared to the Jaccard score.

The domain based scoring method can be used to detect candidate over-clustered families in conjunction with the tree based precision scores. This scoring method can also help in detecting rapidly evolving families.

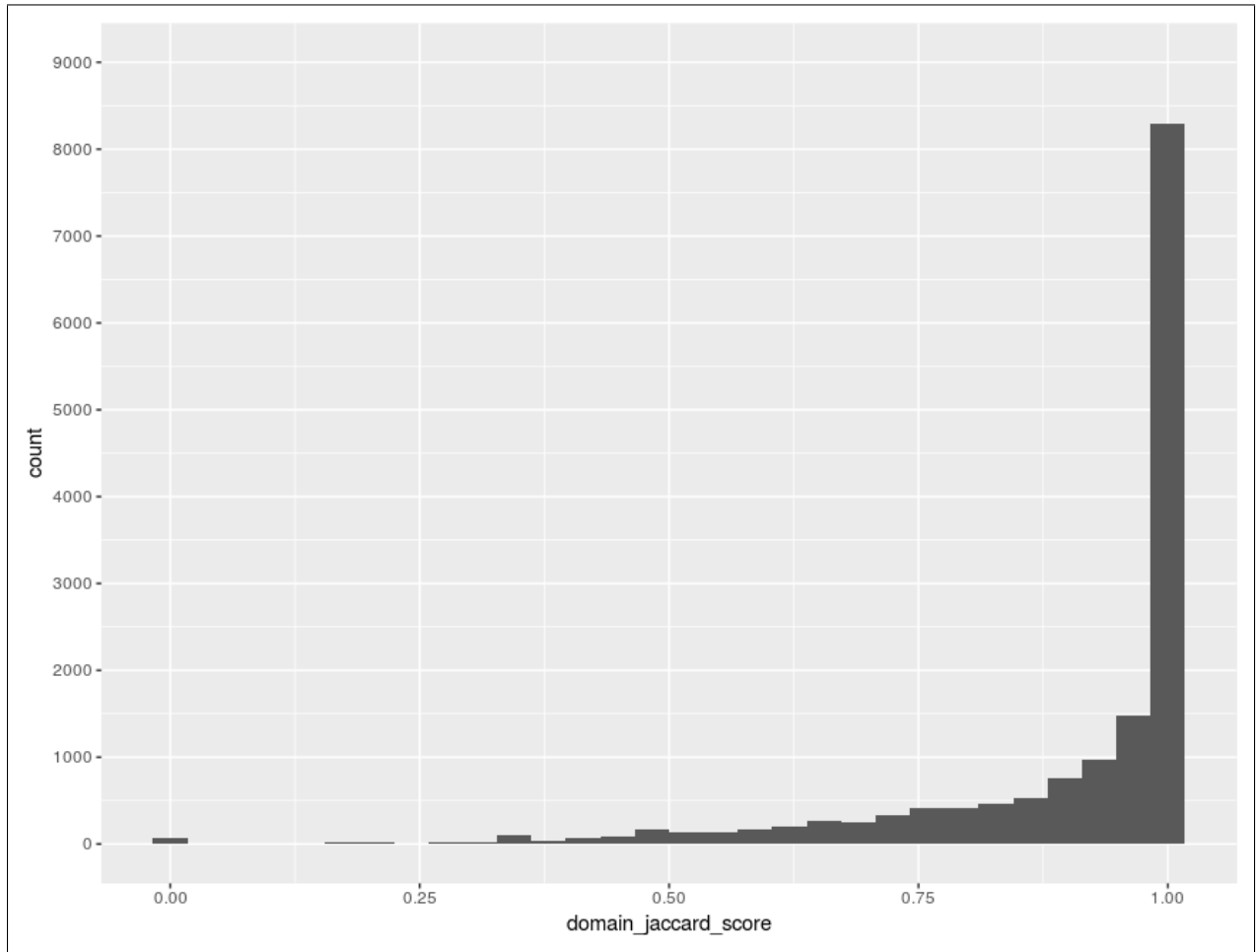


Figure 14: Distribution of domain composition based Jaccard scores for the Ks-based legume families

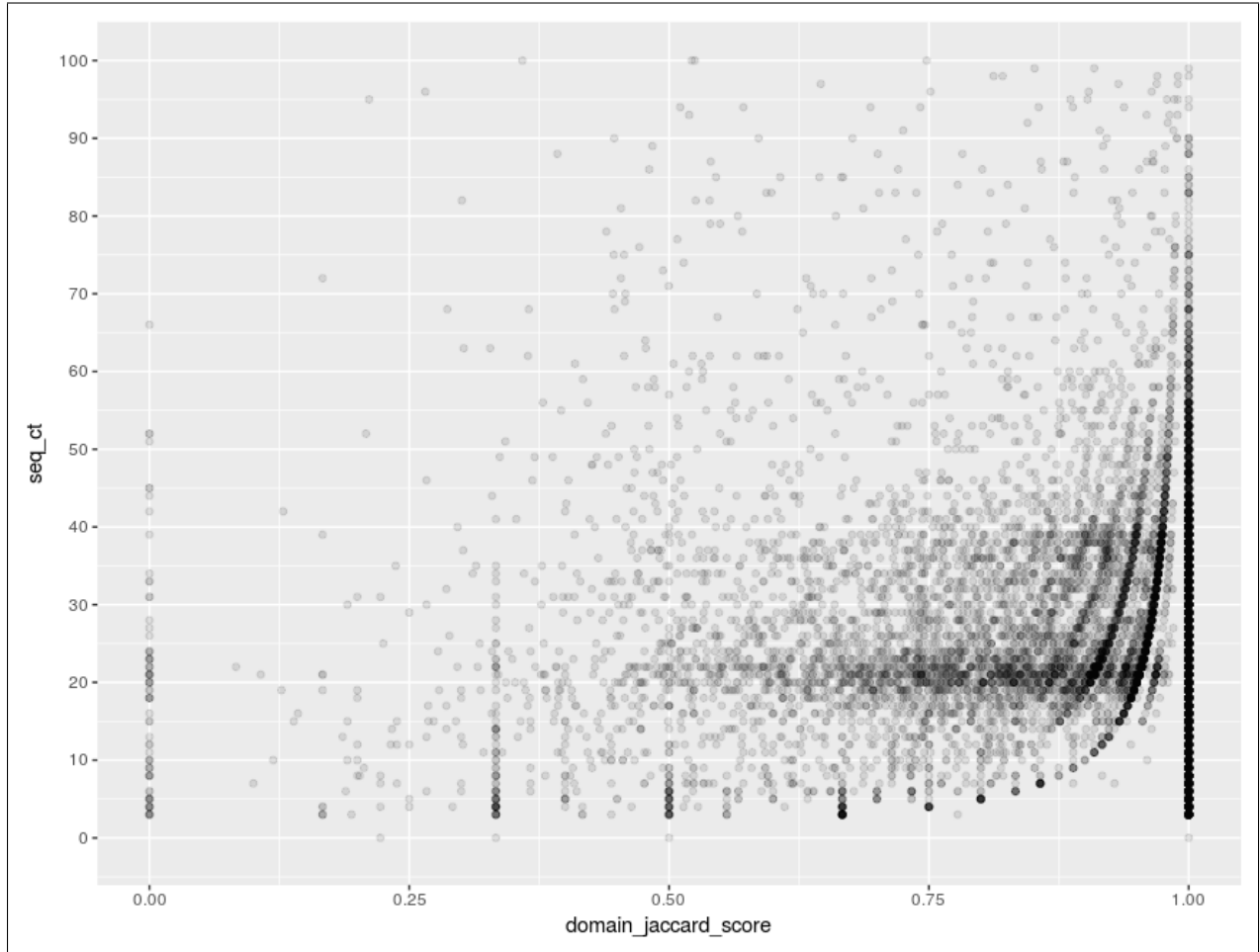


Figure 15: Plot of domain composition based Jaccard score versus ingroup sequence counts for Ks-based legume families

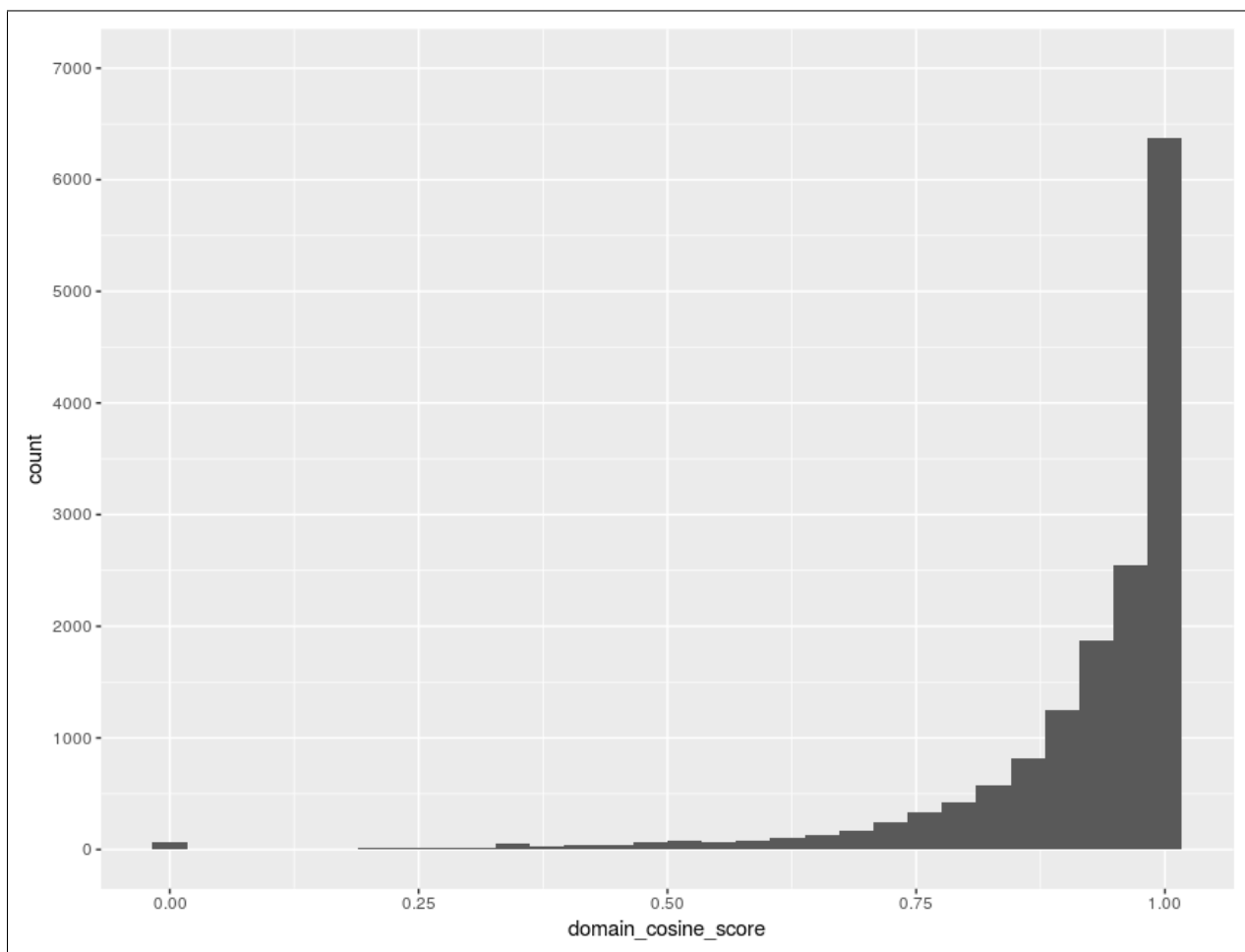


Figure 16: Distribution of domain composition based cosine scores for the Ks-based families

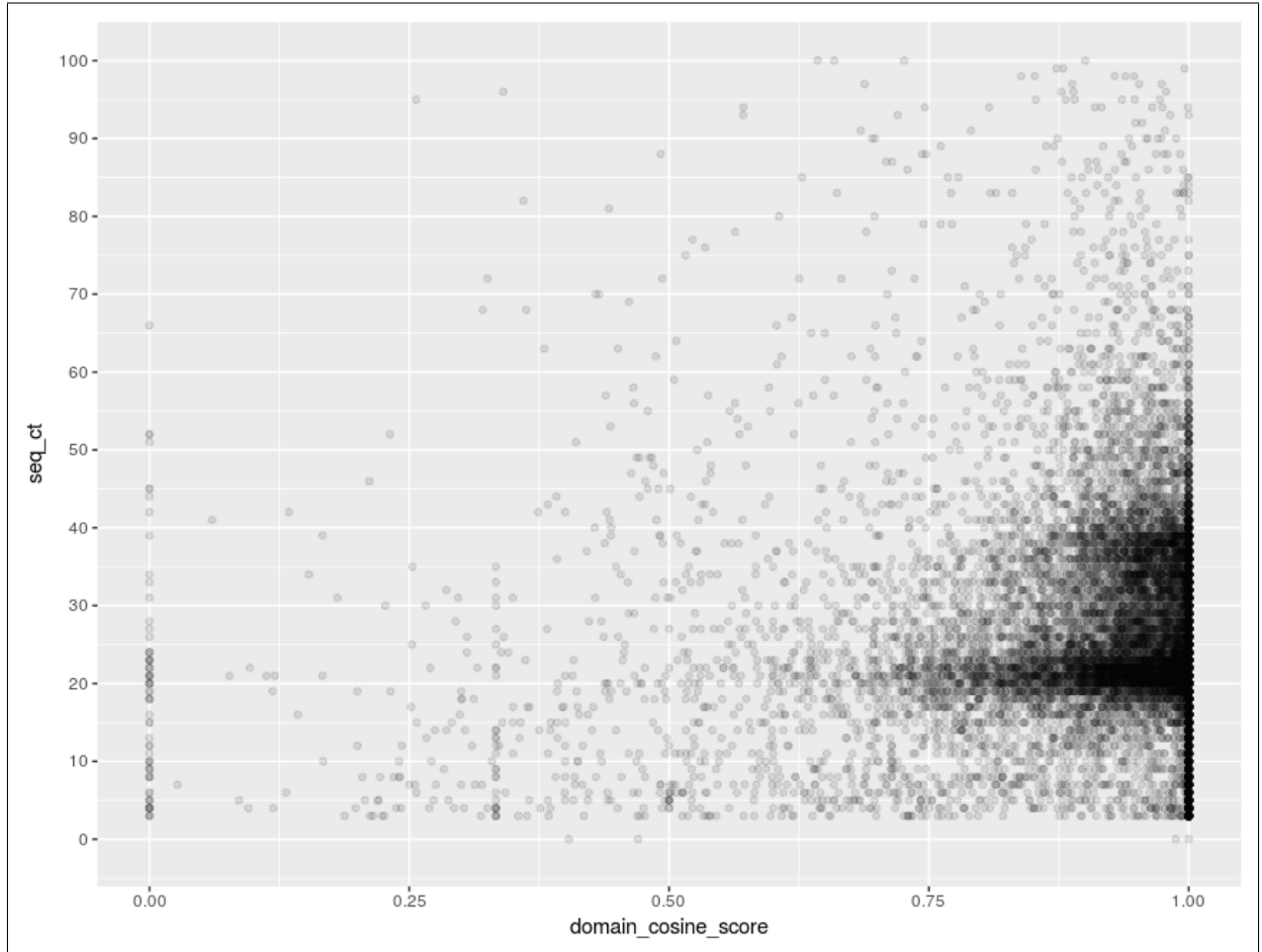


Figure 17: Plot of domain composition based cosine scores versus ingroup sequence counts for Ks-based legume families

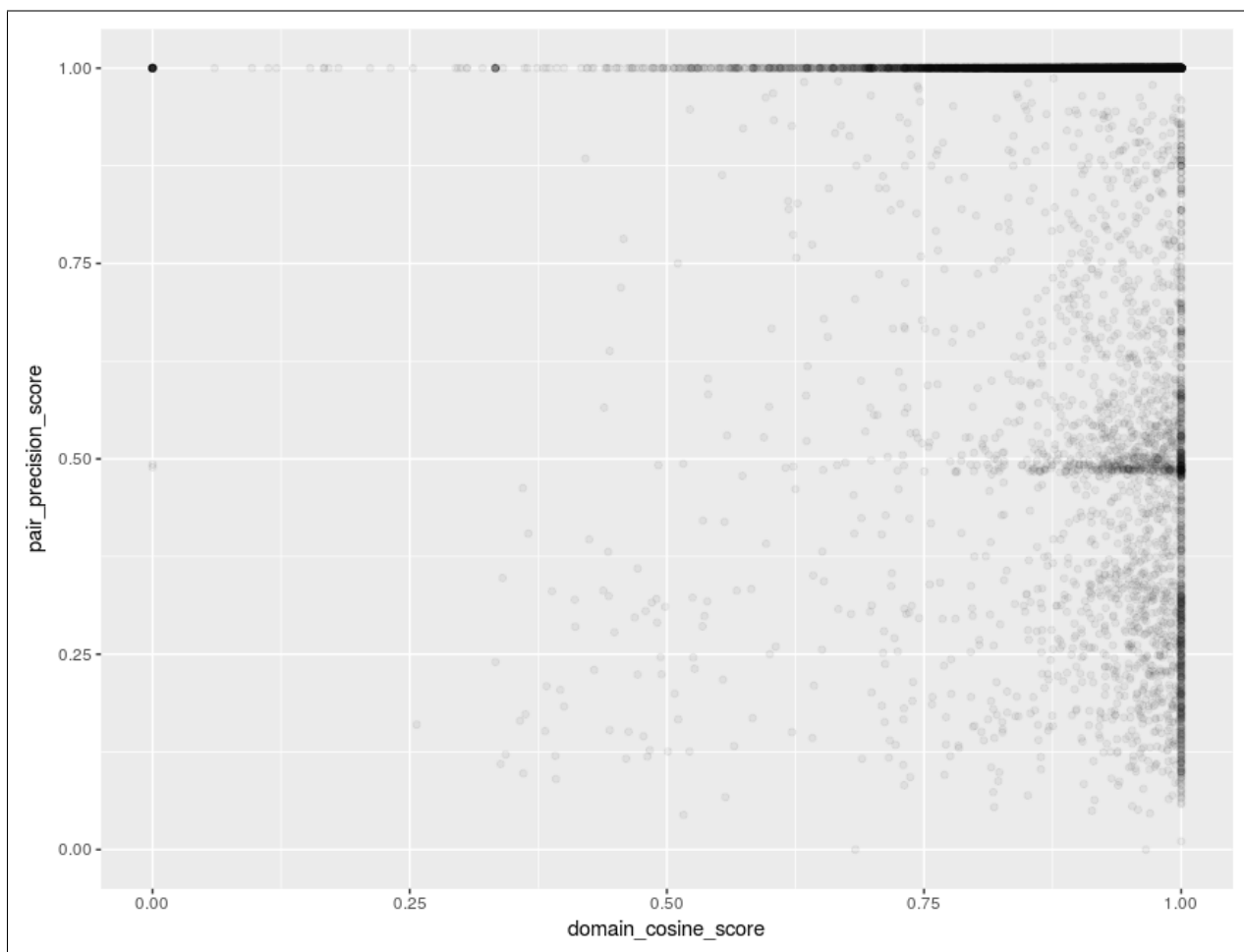


Figure 18: Plot of domain composition based cosine score versus tree based precision scores for Ks-based legume families

### 3.2.4 Future work

Our scoring method will be used in the conjunction with the tree based precision score for detecting and improving legume over-clustered families and for characterizing families. For example, large families with low tree-based precision scores and low domain composition scores are prime candidates for over-clustered families. Also, large families with low tree-based precision scores and high domain composition scores are slowly evolving families that could be over-clustered due to their highly conserved domain architectures over longer evolutionary periods.

It will also be interesting to study biological characteristics of families with high tree-based precision scores and low domain composition scores as these are the rapidly evolving families that have gained or lost protein domains within short evolutionary periods.

## 3.3 Species composition based family scoring

Species composition of a gene family can be defined as the number of sequences represented by per species in the family. Consistent/good families can be expected to have fairly even species compositions with most of the species represented in the family. Also, for well-behaved families the observed sequence counts, per species should be close to expected sequence counts per species. The expected sequence counts for each species per family is largely dependent on occurrences of large-scale duplications and deletions such as WGDs during species evolution.

### 3.3.1 Methods

For a given set of families, built from a specified set of species, a vector of expected species counts per family is obtained using expert opinion. This expected species count vector is compared to the observed species count vector from every family, using cosine similarity function (Equation 1), to assign species composition based scores to corresponding family. Unusually small or large families with highly uneven species compositions are expected to have lower scores than families where most of the species are represented evenly.

$$\text{cosine\_similarity}(O, E) = \frac{O \cdot E}{\|O\| \|E\|} = \frac{\sum_{i=1}^n O_i E_i}{\sqrt{\sum_{i=1}^n O_i^2} \sqrt{\sum_{i=1}^n E_i^2}} \quad (1)$$

### 3.3.2 Preliminary results

The species composition based scoring was applied to Ks-based legume families built from 14 ingroup legume species and 4 outgroup species. The expected count vector corresponding to the 14 legume species was inferred using the histories of the WGD detected in the evolution of the legume species. A WGD has occurred the common ancestor the legume species about 55 Mya (PGWD). Also, lineage specific WGD and WGT have occurred in soybean and lupin, respectively. Moreover, the tetraploid peanut, *Arachis hypogea* is an allopolyploid evolved for the two diploid *Arachis* species. Therefore, 4 soybean genes, 6 Lupin genes, 4 *Arachis hypogea* and 2 genes from rest of the 11 species can be expected to be present in families, assuming no gene loss during family evolution.

The expected count vector was obtained as number of expected number of genes per species (Table 1). All the legume families were scored by comparing the expected counts vector to the observed species counts using cosine similarity function. The plot (Figure 19) of species composition based family scores versus the ingroup sequence counts was obtained to study the relationship between family size that the scores. The plot shows that small families with highly uneven species composition have lower scores and the scores increase as the ingroup sequence count increases.

The legume families fall into at least 2 distinct size classes viz., families which seem to have lost the duplicate gene resulting from the Papilionoid WGD and those which seem to have retained the duplicate gene from the WGD. Families with ingroup sequence counts between 15-45 have highest composition scores with two distinct peaks between 15-30 and 30-45. These peaks correspond to the two types of legume families that have either lost or retained the PWGD duplicate.

Minority of the largest families have high cosine scores. These families have species counts that are far greater than expected, however, the composition is even across all the species. These families are likely result

Species ID	expected gene counts
glyma	4
phavu	2
vigan	2
vigra	2
vigun	2
cajca	2
medtr	2
cicar	2
tripr	2
lotja	2
lupan	6
aradu	2
araip	2
arahy	4

Table 1: Expected gene counts per species per family assuming no gene loss following genome multiplication

of gene-specific duplications that have occurred in the common ancestor. This scoring method rightly assigns higher scores to these families.

### 3.3.3 Significance

The species composition based family scoring assumes that evolutionarily well-behaved families are expected to have sequence counts per species closer to the expected sequence counts. The results on the legume families show that the scoring method assigns higher scores to legume families with sizes between 15-45 confirming the hypothesis that families between the range have even species compositions.

### 3.3.4 Future work

This scoring method can be used in conjunction with pair classification based, domain based and tree based scoring methods for correcting under-clustered and over-clustered legume families. Specifically, small/moderate families with low species composition score can be selected as under-clustered families for further processing through the pair classification approach.

This scoring method can also be used for selecting legitimate ingroup clades while separating/cutting out family clades from over-clustered families detected using the tree based scoring method.



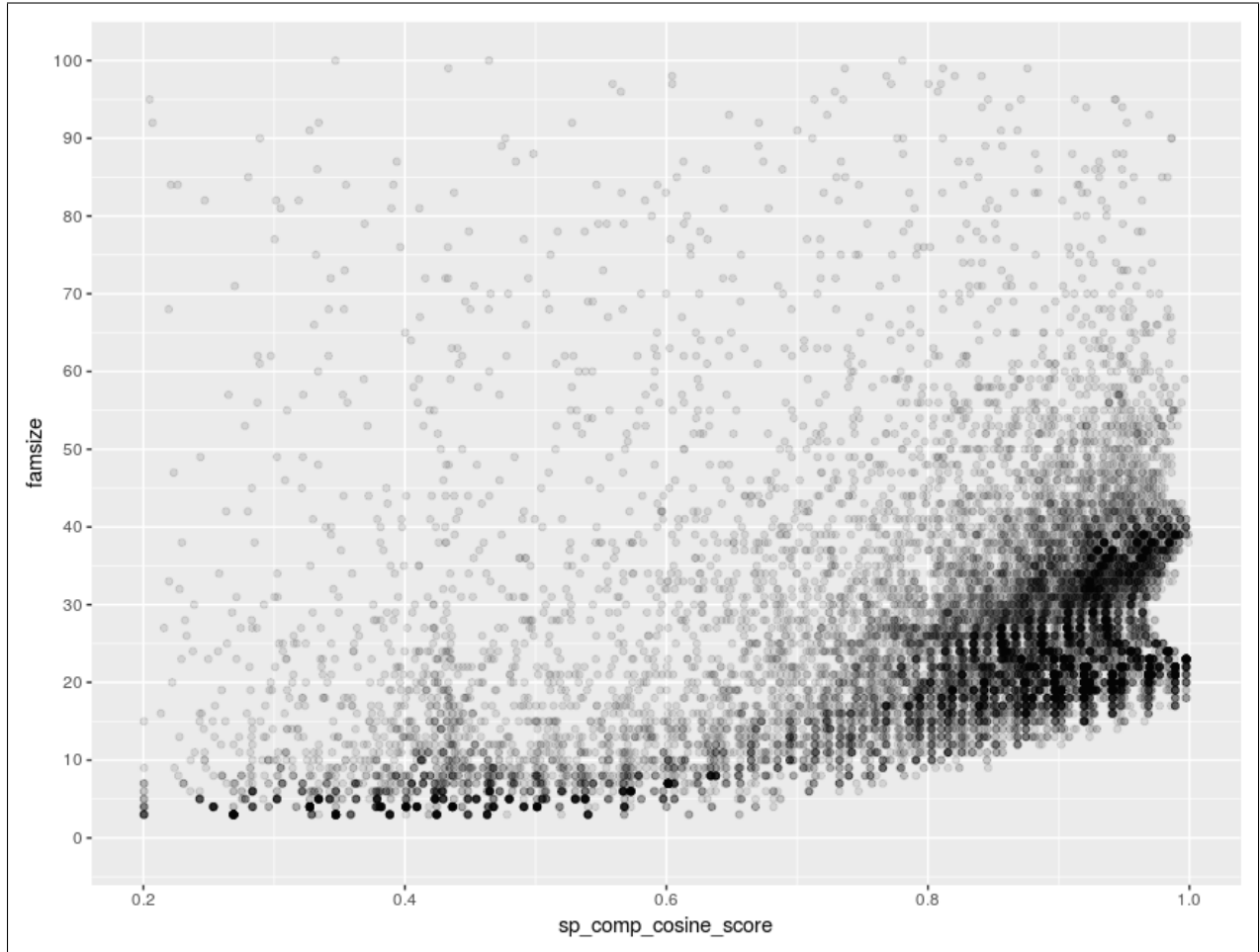


Figure 19: Plot of species composition based family scores versus the ingroup sequence counts of Ks-based legume families

## 4 Objective 3: Develop a workflow for building consensus families from multiple family sets

Family consistency can be assessed by comparing families generated by different family construction methods. This workflow involves comparing two sets of families, obtained from two different family building approaches, to establish family correspondences between the two sets - that is to determine which family from one set best corresponds to which family from the second set. Degrees of overlaps between corresponding families, between the two sets, can be used as measure of family consistency.

The workflow is used to compare set of Ks-based legume families to a set of families obtained from OrthoFinder, to obtain a set of high-quality consensus families. OrthoFinder is one of the latest family building methods that is proven to beat popular family construction methods including OrthoMCL. The two sets are compared through their respective family HMM sets.

### 4.1 Methods

#### 4.1.1 Building HMM sets for families

Gene families are built from 11 legume proteomes with 2 outgroup proteomes, using the OrthoFinder algorithm. Similarly, legume families are also built using in-house Ks-based family building algorithm, from 14 legume proteomes with 4 outgroup proteomes. Two family HMM sets are obtained from both the sets of families using the hmmbuild program from the HMMER package.

#### 4.1.2 Classifying common set of proteomes into families using HMM sets

A common set of 11 legumes proteomes with 2 outgroup proteomes are classified into both the family HMM sets using best matching family HMM strategy. Both HMM sets are searched against the common set of proteomes, using the hmmscan program from HMMER package. Each sequence is classified into any one family, from both the sets, for which the sequence aligns best (e-value  $\leq 1e-5$ ) with the corresponding family HMM.

#### 4.1.3 Calculating family correspondence and family overlaps between the sets

For each family in both the sets, the largest overlapping family in the opposite set is obtained. If, for example,  $fam_1$  from  $set_1$  is the largest overlapping family for  $fam_2$  from  $set_2$  and vice-versa, then  $fam_1$  and  $fam_2$  are considered corresponding families from both the sets. The two overlap scores,  $fam_1 - fam_2$  overlap and  $fam_2 - fam_1$  overlap are also calculated where  $fam_1 - fam_2$  overlap score is the proportion of sequences in  $fam_1$  that overlap with  $fam_2$  and, similarly,  $fam_2 - fam_1$  overlap score is the proportion of sequences in  $fam_2$  that overlap with  $fam_1$ . If both the scores are 1.0 for a pair of families, then the two families match exactly from the two sets.

### 4.2 Preliminary results

The Ks-based and OrthoFinder legume family sets were compared using the family comparison workflow to explore the family correspondence between the two sets.

There were 9160 families that were found to match exactly between both the sets (both overlap scores = 1.0). Eleven thousand five hundred fifty three families were found to match with more than 0.9 overlap score between both sets and 12934 families were found to match more that 0.8 overlap score.

Both the overlap scores for each corresponding family pair from both the sets were plotted against each other (Figure 20). The families near the bottom right of the plot, on the y-axis, are those where the Ks-based family lies entirely within the corresponding orthofinder family. Similarly, the families near top left of the plot, on the x-axis, are those families where the orthofinder family lies entirely within the corresponding Ks-based family.

The distribution of ingroup sequence counts (Figure 21) was also obtained from the families where both the overlap scores were  $\geq 0.9$ , to study the family size distribution of families that are consistently detected

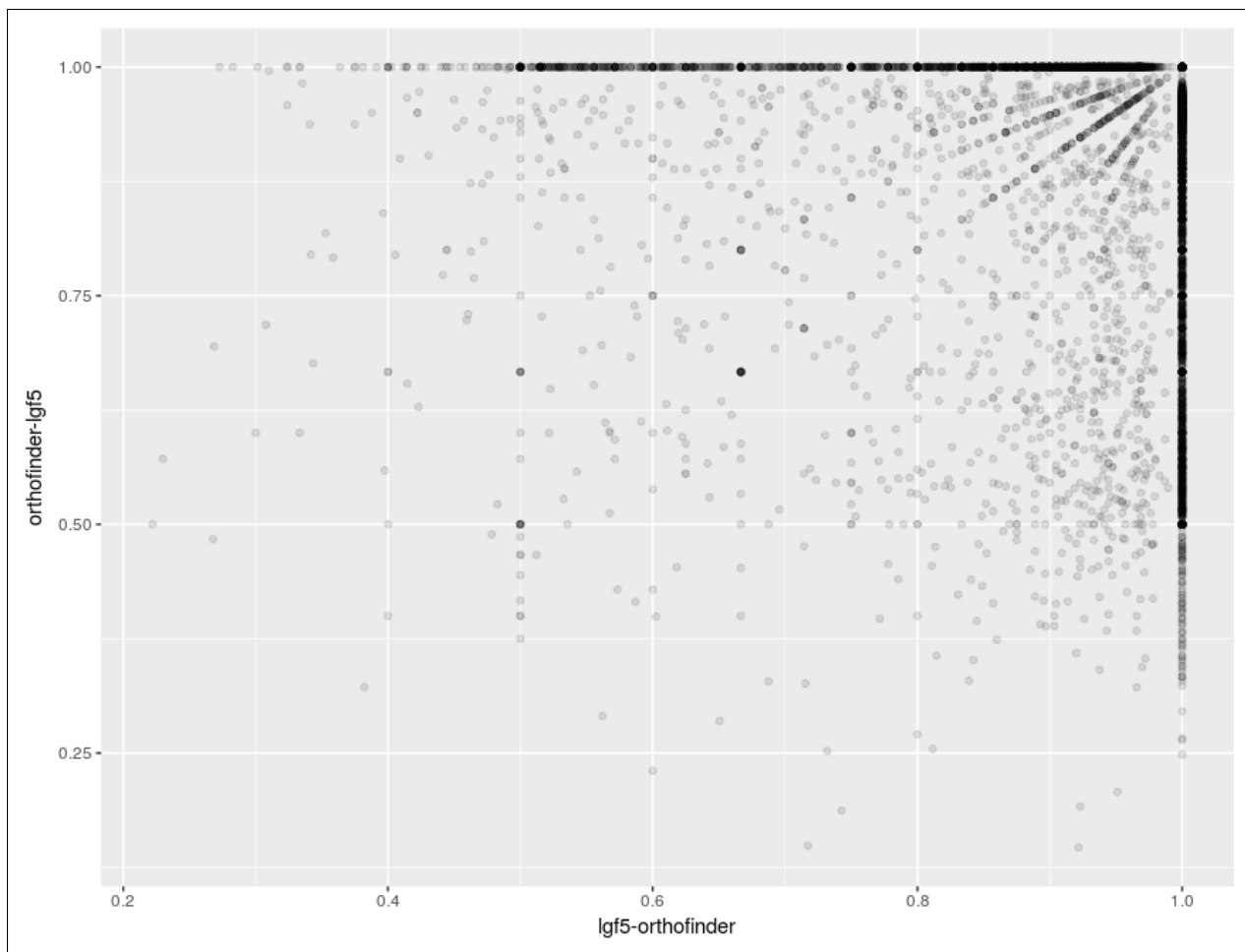


Figure 20: Plot of the two way overlap scores obtained from comparison between Ks-based and OrthoFinder family sets

by both family building algorithms. The distribution shows that consist families are selected from all types of family sizes and not any particular family size.

This workflow has been developed into a PYTHON tool which accepts 2 sets of HMMs corresponding to the 2 family sets and a single set of sequences to be classified into the two family sets.

### 4.3 Future work

The legume families that are found to overlap between the two sets can be released as high-quality set of families. Different sets can be assigned different confidence levels according to the overlaps found between the two sets.

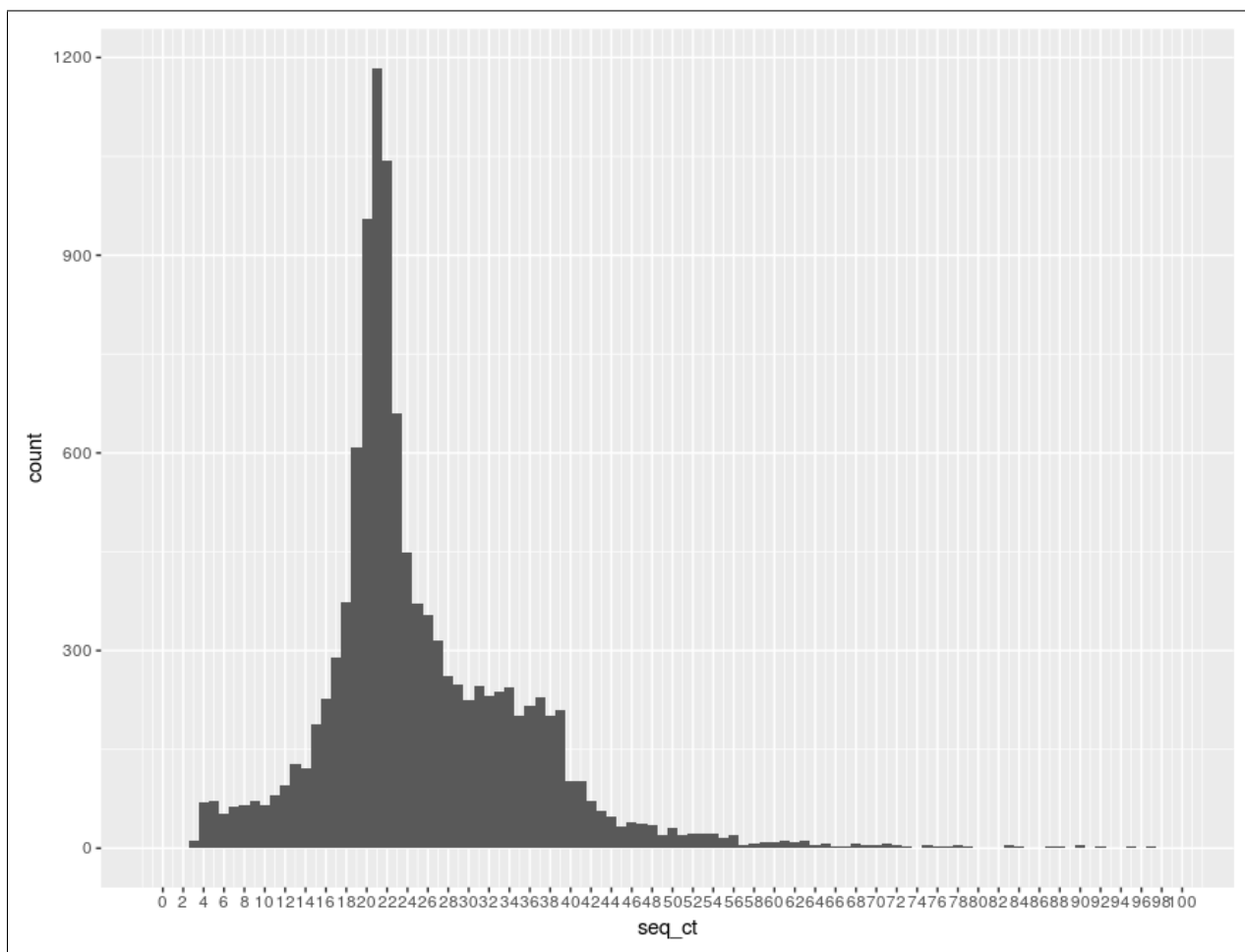


Figure 21: Distribution of in-group sequence counts for Ks-based families with both overlaps scores  $\geq 0.9$

## 5 Objective 4: Biological characterization of the legume families

This section involves exploring biological characteristics of different types of legume families. The legume ancestor is known to have undergone a WGD before divergence (Cannon et al., 2014). However, some families seem to retain the duplicated gene and some families have lost the duplicate gene resulting from the WGD. As a result, legume families can be divided into distinct size classes - families with sizes between 10-30 and families between sizes 30-43. The former class of families are those that seem to have lost the duplicate gene and latter class have retained the duplicate gene. This work investigates if there is any biological basis to this phenomenon, using Gene Ontology (GO) enrichment analysis.

### 5.1 Methods

Legume families are divided into 4 groups corresponding to the different size classes (marked by vertical lines in Figure 22) with group 1 containing families with sizes between 0-10, group 2 containing families with sizes between 10-30, group 3 containing families with sizes between 30-43 and group 4 containing families with sizes 43 and above. Due to exceptionally large number of families in group 2, a subset of families in group 2 is selected for further processing. The dashed lines show the subset of group 2 families.

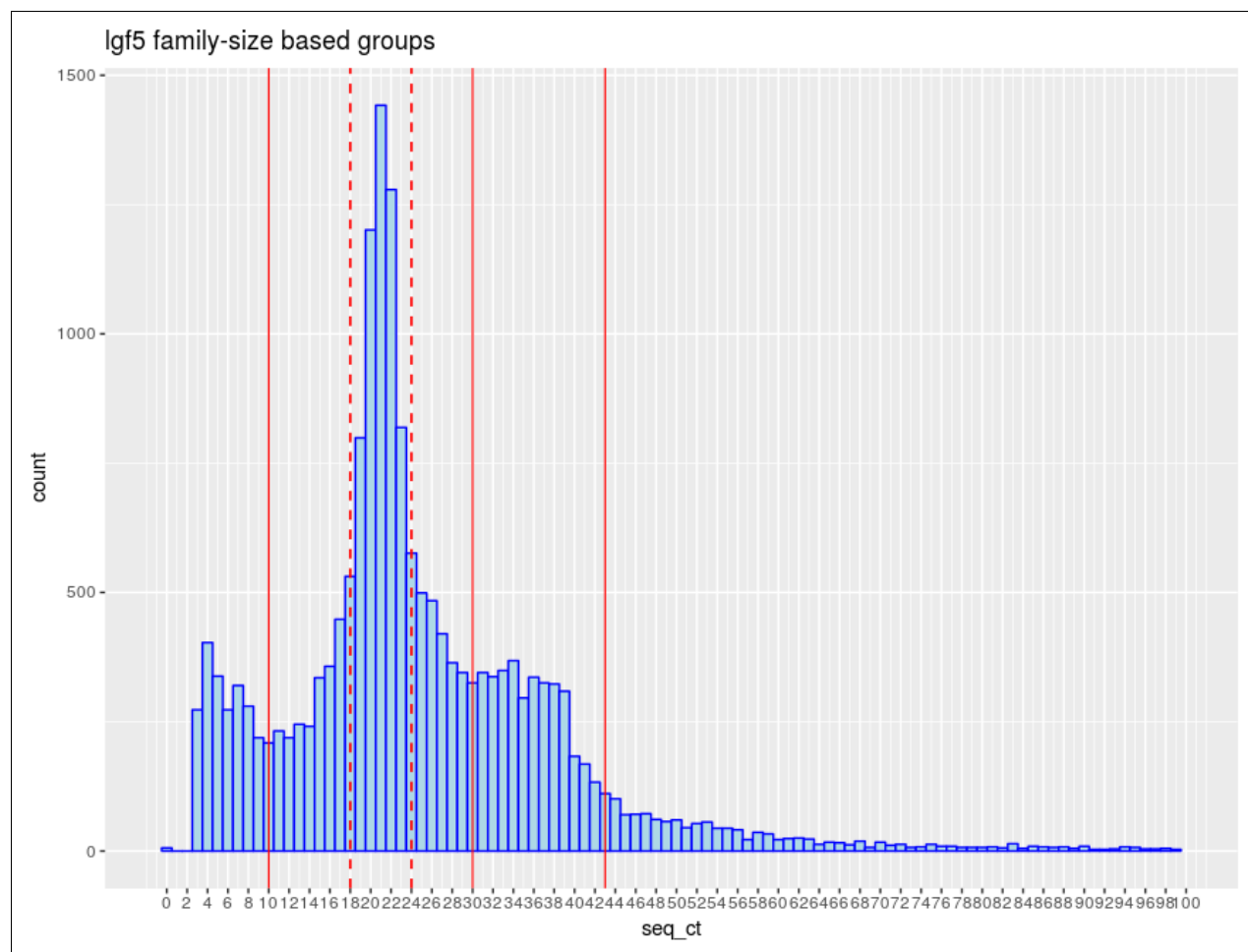


Figure 22: Distribution of in-group sequence counts for Ks-based families

Soybean (*Glycine max*) genes are extracted as representative genes from groups 1, 2-subset, 3 and 4, into separate lists. The gene lists are submitted for GO enrichment analysis on SoyBase (Grant et al., 2009). Enriched ontology terms are obtained for all 4 genes lists corresponding to the 4 family groups.

## 5.2 Preliminary results

Distinct GO terms were found to be significantly enriched for all 4 groups of families. The enrichment for group 1 was less significant. This may be because group 1 may contain fragmented under-clustered families that could be merged into families from other 3 groups. Tables 2, 3, 4, 5 show the top 5 enriched GO terms for each of the size groups.

GO id	P-value	GO description
GO:0009816	$10^{-09}$	defense response to bacterium, incompatible interaction
GO:0010026	$10^{-05}$	trichome differentiation
GO:0006355	$10^{-04}$	regulation of transcription, DNA-dependent
GO:0015690	$10^{-03}$	aluminum cation transport
GO:0009612	$10^{-02}$	response to mechanical stimulus

Table 2: GO enrichment for Group 1 soybean genes

GO id	P-value	GO description
GO:0009658	$10^{-45}$	chloroplast organization
GO:0010027	$10^{-45}$	thylakoid membrane organization
GO:0019288	$10^{-40}$	isopentenyl diphosphate biosynthetic process
GO:0006364	$10^{-37}$	rRNA processing
GO:0019252	$10^{-35}$	starch biosynthetic process

Table 3: GO enrichment for Group 2 subset soybean genes

GO id	P-value	GO description
GO:0006355	$10^{-92}$	regulation of transcription, DNA-dependent
GO:0006468	$10^{-25}$	protein phosphorylation
GO:0010200	$10^{-21}$	response to chitin
GO:0010075	$10^{-20}$	regulation of meristem growth
GO:0009651	$10^{-19}$	response to salt stress

Table 4: GO enrichment for Group 3 soybean genes

GO id	P-value	GO description
GO:0006952	$10^{-191}$	defense response
GO:0007165	$10^{-66}$	signal transduction
GO:0006412	$10^{-54}$	translation
GO:0055114	$10^{-41}$	oxidation-reduction process
GO:0006865	$10^{-40}$	amino acid transport

Table 5: GO enrichment for Group 4 soybean genes

## 5.3 Future work

This analysis can be repeated on the high-confidence legume families obtained from the pair classification based family scoring workflow (merging under-clustered families), tree-based and domain-based scoring workflows (splitting over-clustered families) and the family sets comparison workflow (obtaining highly consistent families).

Applying domain based scoring on the high-confidence families can help in identifying slow-evolving/recent families (high Jaccard and cosine domain composition scores) and fast-evolving/ancient families (low Jaccard

and cosine domain composition scores). GO enrichment analysis of these families can help in understanding the evolutionary pressures experienced by legume species during evolution.

## References

- Alexeyenko, A., Tamas, I., Liu, G., and Sonnhammer, E. L. (2006). Automatic clustering of orthologs and inparalogs shared by multiple proteomes. *Bioinformatics*, 22(14):e9–e15.
- Altenhoff, A. M., Schneider, A., Gonnet, G. H., and Dessimoz, C. (2010). Oma 2011: orthology inference among 1000 complete genomes. *Nucleic acids research*, 39(suppl\_1):D289–D294.
- Bertioli, D. J., Moretzsohn, M. C., Madsen, L. H., Sandal, N., Leal-Bertioli, S. C., Guimarães, P. M., Hougaard, B. K., Fredslund, J., Schauser, L., Nielsen, A. M., et al. (2009). An analysis of synteny of arachis with lotus and medicago sheds new light on the structure, stability and evolution of legume genomes. *BMC genomics*, 10(1):45.
- Bitard-Feildel, T., Kemena, C., Greenwood, J. M., and Bornberg-Bauer, E. (2015). Domain similarity based orthology detection. *BMC bioinformatics*, 16(1):154.
- Blanc, G. and Wolfe, K. H. (2004). Functional divergence of duplicated genes formed by polyploidy during arabidopsis evolution. *The Plant Cell*, 16(7):1679–1691.
- Byrne, K. P. and Wolfe, K. H. (2005). The yeast gene order browser: combining curated homology and syntenic context reveals gene fate in polyploid species. *Genome research*, 15(10):1456–1461.
- Cannon, S. B., McKain, M. R., Harkess, A., Nelson, M. N., Dash, S., Deyholos, M. K., Peng, Y., Joyce, B., Stewart Jr, C. N., Rolf, M., et al. (2014). Multiple polyploidy events in the early radiation of nodulating and nonnodulating legumes. *Molecular biology and evolution*, 32(1):193–210.
- Cannon, S. B., Sterck, L., Rombauts, S., Sato, S., Cheung, F., Gouzy, J., Wang, X., Mudge, J., Vasdewani, J., Schiex, T., et al. (2006). Legume genome evolution viewed through the medicago truncatula and lotus japonicus genomes. *Proceedings of the National Academy of Sciences*, 103(40):14959–14964.
- Davis, J. and Goadrich, M. (2006). The relationship between precision-recall and roc curves. In *Proceedings of the 23rd international conference on Machine learning*, pages 233–240. ACM.
- Dehal, P. S. and Boore, J. L. (2006). A phylogenomic gene cluster resource: the phylogenetically inferred groups (phigs) database. *BMC bioinformatics*, 7(1):201.
- Eddy, S. (1992). Hmmer users guide. *Department of Genetics, Washington University School of Medicine*, 2(1).
- Emms, D. M. and Kelly, S. (2015). Orthofinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome biology*, 16(1):157.
- Falquet, L., Pagni, M., Bucher, P., Hulo, N., Sigrist, C. J., Hofmann, K., and Bairoch, A. (2002). The prosite database, its status in 2002. *Nucleic acids research*, 30(1):235–238.
- Finn, R. D., Tate, J., Mistry, J., Coghill, P. C., Sammut, S. J., Hotz, H.-R., Ceric, G., Forslund, K., Eddy, S. R., Sonnhammer, E. L., et al. (2007). The pfam protein families database. *Nucleic acids research*, 36(suppl\_1):D281–D288.
- Fitch, W. M. (1970). Distinguishing homologous from analogous proteins. *Systematic zoology*, 19(2):99–113.
- Fitch, W. M. (2000). Homology: a personal view on some of the problems. *Trends in genetics*, 16(5):227–231.
- Forterre, P. (2002). A hot story from comparative genomics: reverse gyrase is the only hyperthermophile-specific protein. *Trends in Genetics*, 18(5):236–237.

- Galperin, M. Y. and Koonin, E. V. (1999). Searching for drug targets in microbial genomes. *Current opinion in biotechnology*, 10(6):571–578.
- Gardner, M. J., Hall, N., Fung, E., White, O., Berriman, M., Hyman, R. W., Carlton, J. M., Pain, A., Nelson, K. E., Bowman, S., et al. (2002). Genome sequence of the human malaria parasite plasmodium falciparum. *Nature*, 419(6906):498.
- Grant, D., Nelson, R. T., Cannon, S. B., and Shoemaker, R. C. (2009). Soybase, the usda-ars soybean genetics and genomics database. *Nucleic acids research*, 38(suppl\_1):D843–D846.
- Hubbard, T. J., Aken, B. L., Beal, K., Ballester, B., Caccamo, M., Chen, Y., Clarke, L., Coates, G., Cunningham, F., Cutts, T., et al. (2006). Ensembl 2007. *Nucleic acids research*, 35(suppl\_1):D610–D617.
- Kissinger, J. C., Brunk, B. P., Crabtree, J., Fraunholz, M. J., Gajria, B., Milgram, A. J., Pearson, D. S., Schug, J., Bahl, A., Diskin, S. J., et al. (2002). The plasmodium genome database. *Nature*, 419(6906):490.
- Kroc, M., Koczyk, G., Świecicki, W., Kilian, A., and Nelson, M. N. (2014). New evidence of ancestral polyploidy in the genistoid legume lupinus angustifolius l. (narrow-leaved lupin). *Theoretical and Applied Genetics*, 127(5):1237–1249.
- Letunic, I., Doerks, T., and Bork, P. (2011). Smart 7: recent updates to the protein domain annotation resource. *Nucleic acids research*, 40(D1):D302–D305.
- Lewis, G. P. (2005). *Legumes of the World*. Royal Botanic Gardens Kew.
- Li, H., Coghlan, A., Ruan, J., Coin, L. J., Heriche, J.-K., Osmotherly, L., Li, R., Liu, T., Zhang, Z., Bolund, L., et al. (2006). Treefam: a curated database of phylogenetic trees of animal gene families. *Nucleic acids research*, 34(suppl\_1):D572–D580.
- Li, L., Stoeckert, C. J., and Roos, D. S. (2003). Orthomcl: identification of ortholog groups for eukaryotic genomes. *Genome research*, 13(9):2178–2189.
- Marchler-Bauer, A., Lu, S., Anderson, J. B., Chitsaz, F., Derbyshire, M. K., DeWeese-Scott, C., Fong, J. H., Geer, L. Y., Geer, R. C., Gonzales, N. R., et al. (2010). Cdd: a conserved domain database for the functional annotation of proteins. *Nucleic acids research*, 39(suppl\_1):D225–D229.
- Mistry, J., Bateman, A., and Finn, R. D. (2007). Predicting active site residue annotations in the pfam database. *BMC bioinformatics*, 8(1):298.
- Natale, D., Galperin, M., Tatusov, R., and Koonin, E. (2000a). Using the cog database to improve gene recognition in complete genomes. *Genetica*, 108(1):9–17.
- Natale, D. A., Shankavaram, U. T., Galperin, M. Y., Wolf, Y. I., Aravind, L., and Koonin, E. V. (2000b). Towards understanding the first genome sequence of a crenarchaeon by genome annotation using clusters of orthologous groups of proteins (cogs). *Genome biology*, 1(5):research0009–1.
- Pfeil, B., Schlueter, J., Shoemaker, R., and Doyle, J. (2005). Placing paleopolyploidy in relation to taxon divergence: a phylogenetic analysis in legumes using 39 gene families. *Systematic Biology*, 54(3):441–454.
- Remm, M., Storm, C. E., and Sonnhammer, E. L. (2001). Automatic clustering of orthologs and in-paralogs from pairwise species comparisons. *Journal of molecular biology*, 314(5):1041–1052.
- Roth, A. C., Gonnet, G. H., and Dessimoz, C. (2008). Algorithm of oma for large-scale orthology inference. *BMC bioinformatics*, 9(1):518.
- Schlueter, J. A., Dixon, P., Granger, C., Grant, D., Clark, L., Doyle, J. J., and Shoemaker, R. C. (2004). Mining est databases to resolve evolutionary events in major crop species. *Genome*, 47(5):868–876.
- Söding, J. (2004). Protein homology detection by hmm–hmm comparison. *Bioinformatics*, 21(7):951–960.



- Song, N., Sedgewick, R. D., and Durand, D. (2007). Domain architecture comparison for multidomain homology identification. *Journal of Computational Biology*, 14(4):496–516.
- Sonnhammer, E. L. and Koonin, E. V. (2002). Orthology, paralogy and proposed classification for paralog subtypes. *TRENDS in Genetics*, 18(12):619–620.
- Tatusov, R. L., Galperin, M. Y., Natale, D. A., and Koonin, E. V. (2000). The cog database: a tool for genome-scale analysis of protein functions and evolution. *Nucleic acids research*, 28(1):33–36.
- Tatusov, R. L., Koonin, E. V., and Lipman, D. J. (1997). A genomic perspective on protein families. *Science*, 278(5338):631–637.
- Tatusov, R. L., Natale, D. A., Garkavtsev, I. V., Tatusova, T. A., Shankavaram, U. T., Rao, B. S., Kiryutin, B., Galperin, M. Y., Fedorova, N. D., and Koonin, E. V. (2001). The cog database: new developments in phylogenetic classification of proteins from complete genomes. *Nucleic acids research*, 29(1):22–28.
- Van der Heijden, R. T., Snel, B., Van Noort, V., and Huynen, M. A. (2007). Orthology prediction at scalable resolution by phylogenetic tree analysis. *BMC bioinformatics*, 8(1):83.
- Wapinski, I., Pfeffer, A., Friedman, N., and Regev, A. (2007). Automatic genome-wide reconstruction of phylogenetic gene trees. *Bioinformatics*, 23(13):i549–i558.
- Whetzel, P. L., Date, S. V., Gajria, K., Fraunholz, M. J., Gajria, B., Grant, G. R., Iodice, J., Labo, P. T., Milgram, A. J., Stoeckert, C. J., et al. (2005). Plasmodb: The plasmodium genome resource. In *Molecular Approaches to Malaria*, pages 12–23. American Society of Microbiology.