

## Abstract

Given a set of examples belonging to different classes, the goal of supervised taxonomy generation is to identify class-uniform groups of organisms with close phylogenetic relatedness that corresponds to subclasses of the underlying class structure. By extracting such sub-trees from supervised taxonomies, prediction models that infer class membership of unlabeled organisms can ultimately be developed. Additionally, by analyzing the generated tree structure, a biologist can interpret the clustering result and gain insights into how the biological groupings are related.

## Problem Statement and Research Goals

- Traditional phylogenetic-based taxonomy generation approaches rely on gene sequence proximities to merge clusters into larger clusters
- Supervised Taxonomies (S.T.)** on the other hand, are generated using both sequence proximity and background knowledge in the form of class labels.
- Given organisms belonging to different classes, supervised taxonomy aims to identify class-uniform groups that correspond to subclasses of the underlying class structure.

(a) Dataset

	<i>D. melanogaster</i>	0	0.138	0.199	0.248	0.238	0.211	0.245	0.245	0.234	0.252	0.251	0
<i>D. pseudoobscura</i>	0.138	0	0.199	0.24	0.241	0.197	0.235	0.232	0.227	0.243	0.227	0	0
<i>S. lebanonensis</i>	0.199	0.199	0	0.256	0.259	0.211	0.244	0.247	0.239	0.245	0.234	0	0
<i>S. albivittata</i>	0.248	0.24	0.256	0	0.092	0.213	0.18	0.176	0.167	0.178	0.173	0	0
<i>D. crassifemur</i>	0.238	0.241	0.259	0.092	0	0.214	0.182	0.192	0.185	0.19	0.182	0	0
<i>D. mulleri</i>	0.211	0.197	0.211	0.213	0.214	0	0.177	0.177	0.177	0.18	0.18	0	0
<i>D. affinisdisjuncta</i>	0.245	0.235	0.244	0.18	0.182	0.177	0	0.037	0.06	0.06	0.091	0	0
<i>D. heteroneura</i>	0.245	0.232	0.247	0.176	0.192	0.177	0.037	0	0.052	0.059	0.083	0	0
<i>D. mimica</i>	0.234	0.227	0.239	0.167	0.185	0.177	0.06	0.052	0	0.056	0.075	0	0
<i>D. adiasstola</i>	0.252	0.243	0.245	0.178	0.19	0.18	0.06	0.059	0.056	0	0.091	0	0
<i>D. nigra</i>	0.251	0.227	0.234	0.173	0.182	0.18	0.091	0.083	0.075	0.091	0	0	0

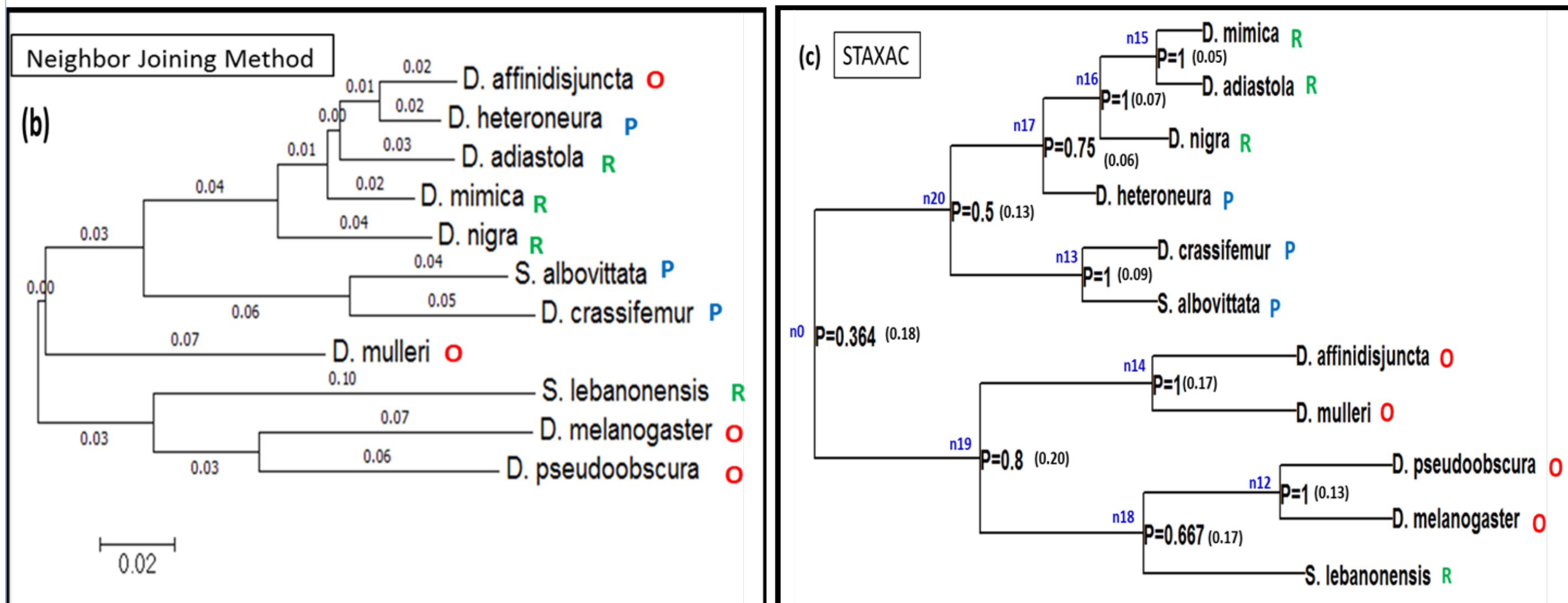


Fig. 1: Example of Tree Generated by S.T. Vs Tree Generated by Traditional Method (N.J)

## Algorithms

Two Supervised Taxonomies (S.T.) generation algorithms are proposed:  
Input: Dissimilarity matrix (*pairwise distances among the taxa*)  
Output: Binary tree

- Bottom up approach called STAXAC (**S**upervised **T**axonomy **A**gglomerative **C**lustering)
- Top Down approach called MMCSC (**M**edoid of the **M**ajority **C**lass **S**upervised **C**lustering)

### Algorithm 1: STAXAC

Input: examples with class labels and their dissimilarity matrix D.

Output: Hierarchical clustering

- Start with a clustering X of one-object clusters.
- $\forall C, C' \in X$  merge-candidate( $C, C'$ )  $\leftrightarrow$  ( $1 - \text{NN}_X(C') = C$  or  $1 - \text{NN}_X(C) = C'$ )
- WHILE there are merge-candidates ( $C, C^*$ ) left

BEGIN

- Merge the pair of merge-candidates ( $C, C^*$ ) obtaining a new cluster  $C' = C \cup C^*$  and a new clustering  $X'$  for which  $\text{Purity}(X')$  has the largest value
- Update merge-candidates:  
 $\forall C''$  merge-candidate( $C', C''$ )  $\leftrightarrow$  (merge-candidate( $C, C''$ ) or merge-candidate( $C^*, C''$ ))

Extend dendrogram by drawing edges from C and  $C^*$  to  $C'$

END

- Return constructed dendrogram

### Efficient implementation

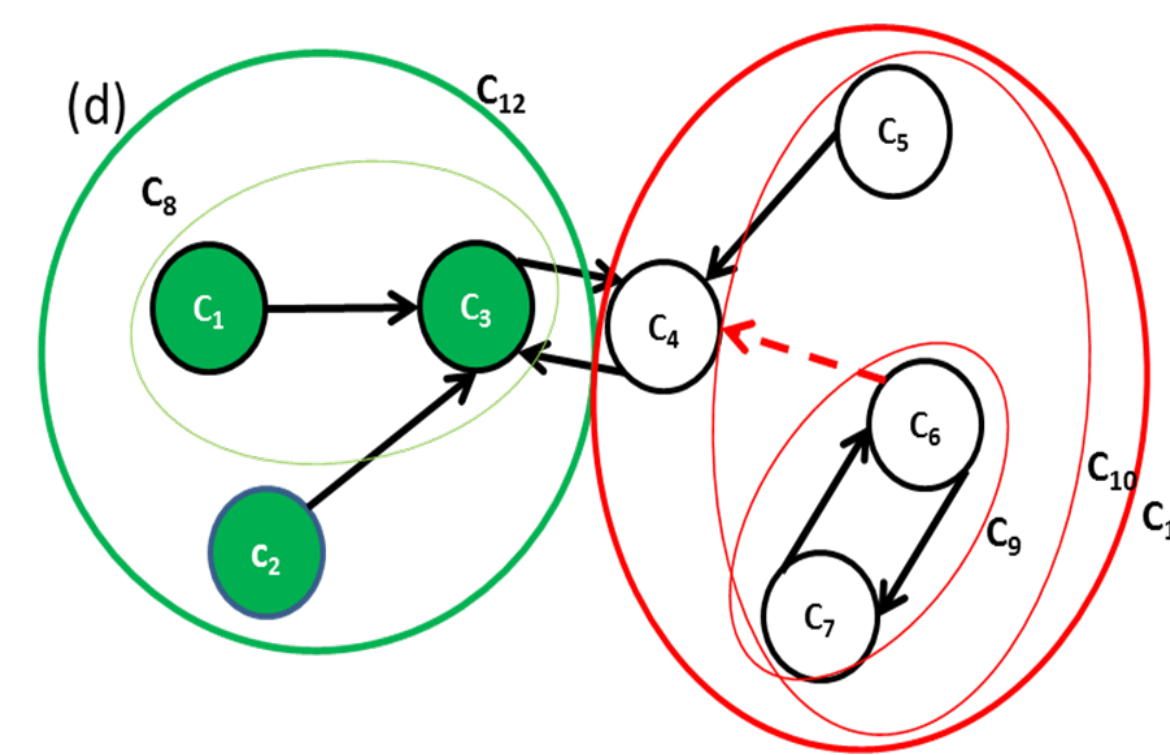


Fig. 2: S.T. Approach Generates High Purity Clusters

Definition 1: Given  $C \in X$ , merge-candidate set of C  
 $\mathcal{M}(C) = \{C' \mid C' \in 1 - \text{NN}_X(C) \vee C' \in 1 - \text{NN}_X^{-1}(C)\}$  (1)

Definition 2: Given  $C \in X, C' \in X$  and  $C^* \in X$  such that  $C = C' \cup C^*$   
 $\mathcal{M}(C' \cup C^*) = \mathcal{M}(C') \cup \mathcal{M}(C^*) - \{C', C^*\}$  (2)

Keep track of  $1\text{NN}_X$  relationships as follows  
After each merge operation  
update the  $\mathcal{M}(C' \cup C^*)$  using equation (2)  
If  $\mathcal{M}(C' \cup C^*) = \emptyset$  then  
 $\mathcal{M}(C' \cup C^*) \leftarrow \text{compute}(1\text{NN}_X(C))$

### Algorithm 2: MMCSC

Inputs: a set of examples with class labels and their dissimilarity matrix D; purity threshold  $\theta$ .

Function: MMCSC(D,  $\theta$ )

IF  $\text{Purity}(D) > \theta$  terminate;

Separate the examples in tree node D into two subsets: one containing only sequences of the majority class  $D_{\text{majority}}$  and the other composed of the remaining sequences,  $D_{\text{not\_majority}}$

Compute the medoids of both sets:  $\text{medoid}(D_{\text{majority}})$ ,  $\text{medoid}(D_{\text{not\_majority}})$

#Split node D in two new nodes:  $D1, D2$

FOR each d in D

IF d is closer to  $\text{medoid}(D_{\text{majority}})$  than it is to  $\text{medoid}(D_{\text{not\_majority}})$

assign d to D1

ELSE assign d to D2

Call MMCSC(D1,  $\theta$ )

Call MMCSC(D2,  $\theta$ )

END Function

## Experimental Results

We evaluate STAXAC, and MMCSC against each other and against NJ (Neighbor Joining) a popular distance-based Phylogenetic approach using two criteria:

□ Runtime

□ Tree Topology

- Average Purity Per Depth

- Number of Sub-family Changes

(From root node to a leaf-node, count number of changes in majority class label)

- Edited Tree Sizes

(Merge leaf-nodes from bottom up until leaf-node purity  $\leq$  threshold  $\theta$ )

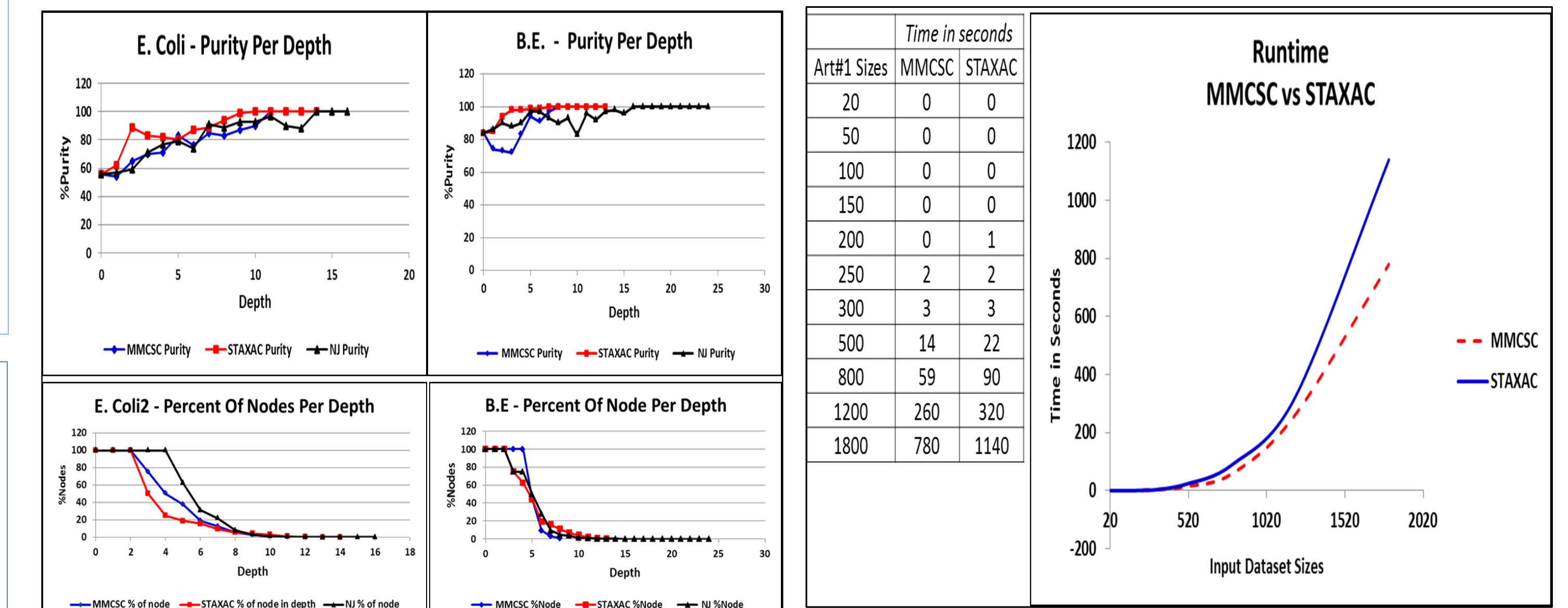


Fig. 3: Average Purity Per Depth

Fig. 4: Runtime

TABLE I: Number Of Sub-family Changes

	MMCS	STAXAC	NJ
E. coli	34	18	35
BE	15	7	21
BV	53	26	56
BH	8	4	8
Art#2	13	4	49
Art#3	36	6	61

TABLE II: Edited Tree Sizes

Dataset name	MMCS			STAXAC			NJ		
	$\theta = >$	100%	95%	100%	95%	85%	100%	95%	85%
E. coli	103	103	101	55	55	48	123	123	123
BE	57	57	51	25	22	11	73	73	71
BV	197	197	197	111	51	51	273	273	269
BH	57	57	57	19	17	3	75	73	65
Art#2	71	71	71	15	15	12	139	139	139
Art#3	109	109	105	17	17	17	153	153	153

## Conclusion & Future Work

Both algorithms can help a biologist interpret the clustering result and gain insights into how the clustering (biological groupings) are related in the context of an underlying class structure.

### Future Work:

- Predictive Models using as input trees generated by S. T. approach
- Using trees generated by S. T. to determine if a phenotype (class label) within a closely related species is genetically induced.