

# Phylogenetic relationships and characterization of UDP-glycosyltransferases (UGTs) in dicotyledonous plants

Hugo Marcelo Atencio; Nicolas Stocchi; Agustín Amalfitano; Fernando Villarreal; Arjen ten Have



I I B



## BACKGROUND

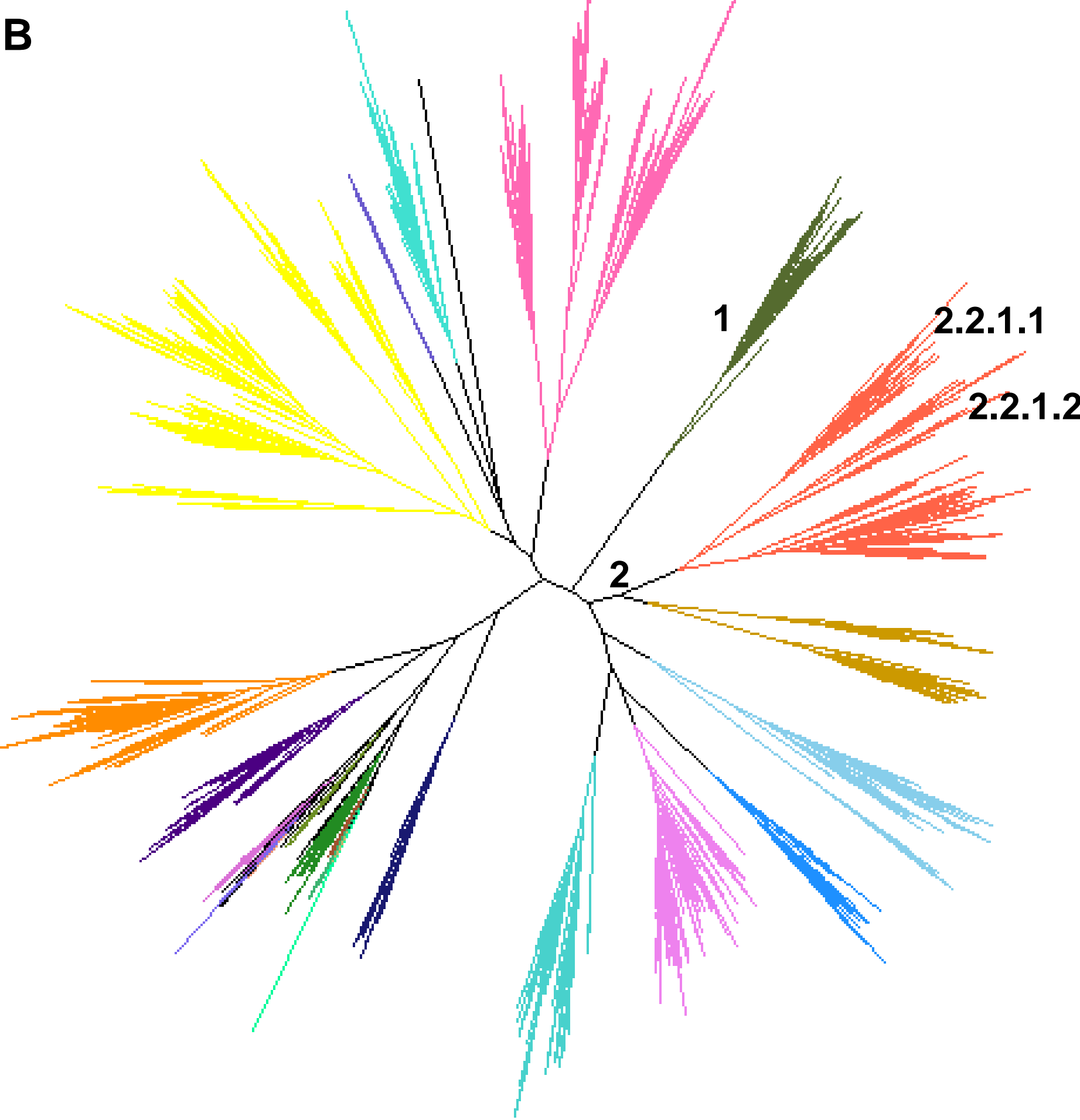
- The protein superfamily UDP-Glycosyl Transferase (UGT) is highly divergent in plants.
- UGTs transfer sugar moieties from a variety of UDP-activated sugars donor to a variety of acceptor molecules.
- Plant UGTs have a consensus motif of 44 amino acids denoted Plant Secondary Product Glycosyltransferase (PSPG motif).
- UGTs are important in secondary metabolism and they have evolved towards a complex superfamily.
- Substrate promiscuity has likely resulted in erroneous function assignation and might explain why many UGT phylogenies do not correspond with functional classification.
- We are interested in plant UGTs involved in flavonoid and anthocyanin synthesis and need to identify the paralogs involved.
- In the initial sequence mining we identified many partial sequences as well as many pseudogenes in datasets.
- The presence of these “bad” sequences in datasets results in erroneous MSA and therewith erroneous phylogenies.

## OBJECTIVES

- > To reconstruct and characterize a high fidelity phylogeny for structure-function prediction of plant dicotyledonous UGTs.
- > To identify Specific Determining Positions (SDPs) of UGT, involved in the catalytic or structural function on anthocyanin metabolism.

## CONCLUSIONS

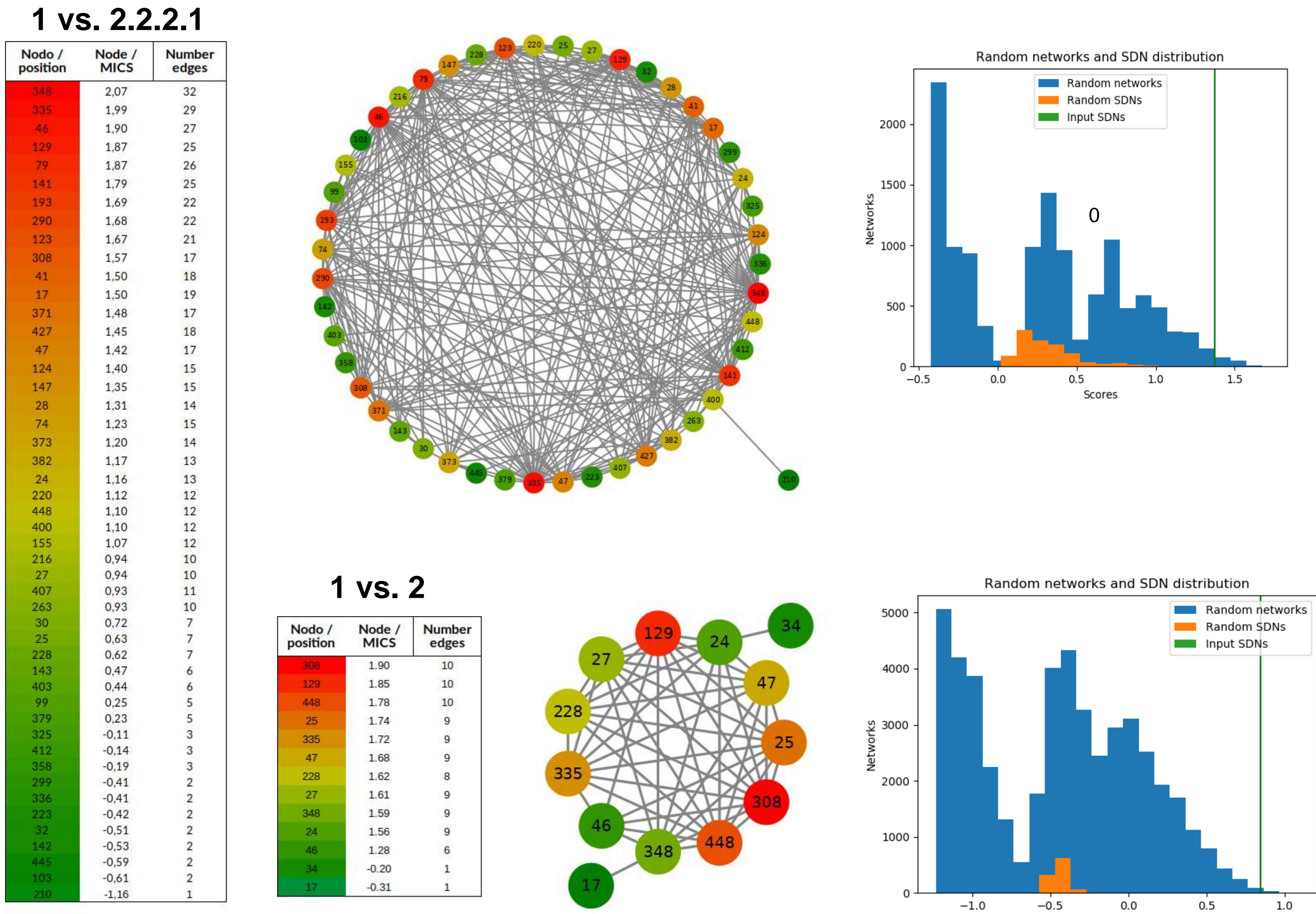
- > A high fidelity phylogenetic tree of UGTs was obtained and will contribute to a more precise functional annotation based on studies of structure-function prediction.
- > The UGTs involved in anthocyanidin and anthocyanin glycosilation (3-O-UGT and 5-O-3UGT respectively, are NOT monophyletic despite the clear substrate similarity.
- > Identification of SDPs in hampered by the fact that the 3-O-UGT subfamily has no sister clade, by which comparison has to be done indirectly.



Seqrutinator / module order	Sequences removed	Sequences accepted
1 <sup>st</sup> SSR	1783	-
2 <sup>nd</sup> NHHR	24	-
3 <sup>rd</sup> GIR	113	-
4 <sup>th</sup> CGSR	632	-
5 <sup>th</sup> PsOR	42	-
Total (4950)	2693	2257

Cluster/ Subcluster	MIT	Network MICS	Network nodes
1 vs. 2	7.7	1.37	13
1 vs. 2.2.1.1	5.32	0.85	48

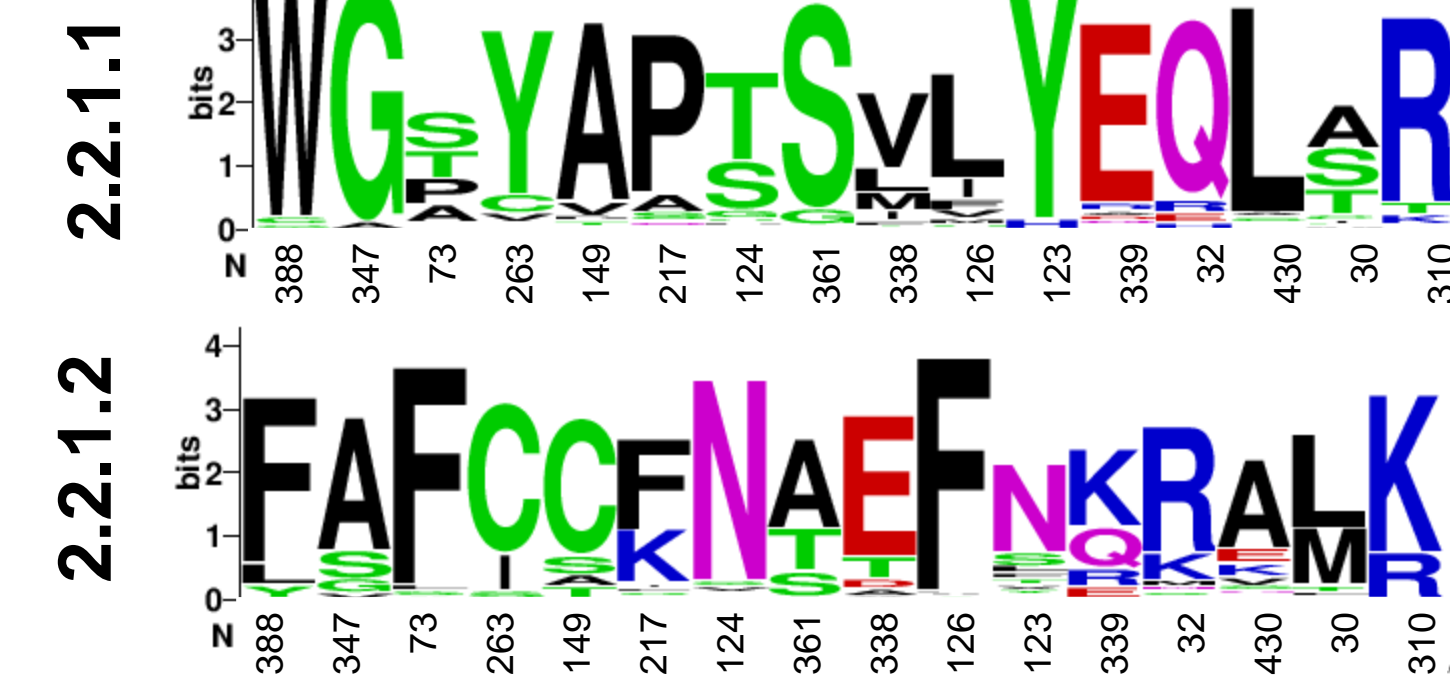
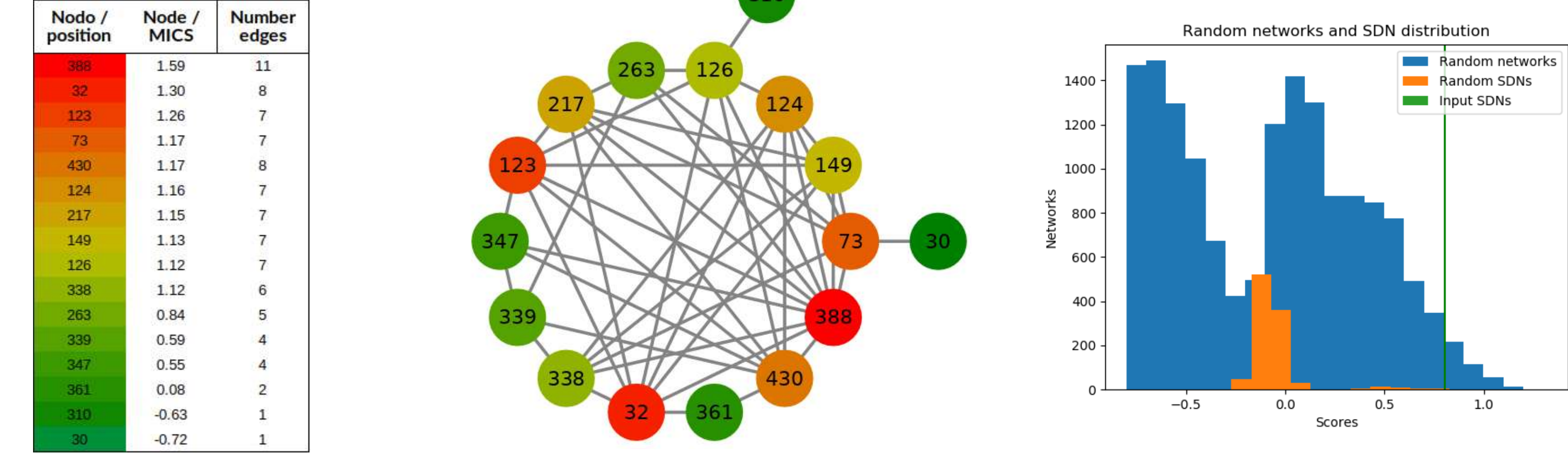
Subcluster	MIT	Network Score	Network nodes
2.2.1.1 vs. 2.2.1.2	4.28	0.8	16



### 1)- UGT phylogenetic tree and clustering for identification of CDPs and SDPs.

A Of 4950 homologs, a total of 2257 UGT sequences passed SEqrutinator (See table for details) and were considered bona fide sequences used for Multiple Sequence Alignment phylogeny. B HMMERCTTER clustering using the obtained high quality phylogeny resulted in 24 clusters that show 100% Precision and Recall. The minor clade 1 corresponds with 3-O-UGT whereas 2.2.1.1 is 5-O-3UGT (UGTs that glycosylate at position 3 and 5 of anthocyanidins and anthocyanin, respectively). Note that these subfamilies evolved from different ancestors, which hampers the identification of SDPs using SDPFox, Mistic and our novel Specificity Networks Mutual Information Connectivity Score (MICS) C.

### 2.2.1.1 vs. 2.2.1.2



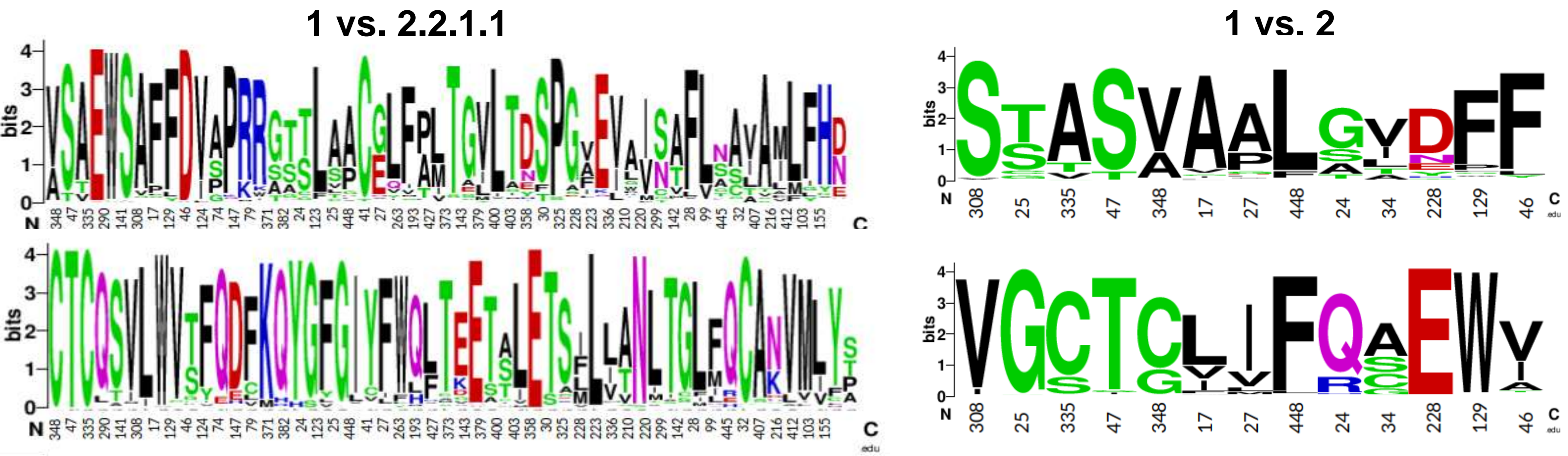
### 4)- Contrast within the monophyletic clade.

Comparing 5-O-3UGT (2.2.1.1) with sister clade (2.2.2.2). Only 4 of the 13 predicted SDPs are shared with SDPs identified between clades 1 vs. 2.2.1.1, but do not match those SDPs detected by comparing 1 vs 2.

### 2)- Structure Function Prediction of Distantly related Clades.

Comparing the 3-O-UGT (1) with 5-O-3UGT (2.2.1.1) is hampered by the lack of monophyly. SDPs identified might be involved in the diversification of clades 1 and 2, rather than the diversification of 3-O-UGT (1) with 5-O-3UGT (2.2.1.1) only. The green line on the two plots “Random networks and SDN distribution” corresponds with the MICS of the SDP network. A total of 48 SDPs are identified when comparing 1 with 2.2.1.1., of which many might actually result from the diversification of 1 and 2. Indeed, the 13 SDPs identified by comparing 1 with 2, are identified when comparing 1 with 2.2.1.1.

Although the SDPs identified in the 1-2 comparison are not exclusive to the diversification we analyzed, they might involve mutations that have allowed for the diversification toward the 5-O-3UGT activity. A more careful and stepwise of the family is required to obtain more reliable predictions.



### 3)- Specificity Determining Network (SDN): selection of sensitive and specific threshold.

Application of the MIT at 5.32 identified a fully connected network has the highest Score SDN= 0.85 and contains 48 SDPs/Nodes between 1 vs 2.2.1.1. The MIT determined for clades 1 vs 2 was 1.37 with the largest network MICS. The logo sequences for each of the comparisons are shown.

## Data Mining and Methods

- UGT homologues with HMMER. 17 complete proteomes of dicotyledonous plants and the control-set of SwissProt. - 2257 UGT homologs sequences scrutinized objectively with SEqrutinator. - MSA: MAFFT GINS-i, Trimming: BMGE. - Phylogeny: PHYML; Clustering: HMMERCTTER (100% P&R), CDPs (Cluster determining Positions): SDPFox. - MI: mutual information (MISTIC: <http://mistic.leloir.org.ar/index.php>). - MICS: mutual information connectivity score. -  $MICS = \log_{10} \left( \frac{C * [\Sigma(1,5^*MI) - MIT]}{N - 1} \right)$ . MI= mutual information, MIT= mutual information threshold, N= number of nodes, C= connectivity. - Network Score: average of MICS (SDN score is compared to 1000 randomly selected, connected subnetworks inferred by the same MIT). - Logos: <https://weblogo.berkeley.edu/logo.cgi>.